

13.2.3 Discrete Time Convolutions

D. (Discrete Convolution)

For $f, h: \mathbb{R}^d \rightarrow \mathbb{R}$, we define the discrete convolution via

$$(f * h)[u] := \sum_{t=0}^{\infty} f[t]h[u-t] = \sum_{t=0}^{\infty} f[u-t]h[t]$$

Com. Note that the use of rectangular brackets suggests that we're using "arrays" (discrete-time samples).

Com. Typically we use h with finite support (windows size).

D. (Multidimensional Discrete Convolution)

For $f, h: \mathbb{R}^d \rightarrow \mathbb{R}$ we have

$$(f * h)[u_1, \dots, u_d] := \sum_{t_1=0}^{\infty} \dots \sum_{t_d=0}^{\infty} f[t_1, \dots, t_d]h(u_1 - t_1, \dots, u_d - t_d)$$

$$= \sum_{t_1=0}^{\infty} \dots \sum_{t_d=0}^{\infty} f[u_1 - t_1, \dots, u_d - t_d]h(t_1, \dots, t_d)$$

D. (Discrete Cross-Correlation)

Let $f, h: \mathbb{R}^d \rightarrow \mathbb{R}$

$$(h * f)[u] := \sum_{t=0}^{\infty} h[t]f[u+t] = \sum_{t=0}^{\infty} h[-t]f[u-t]$$

$$= (\bar{h} * f)[u] = (f * \bar{h})[u] \quad \text{where } \bar{h}(t) = h(-t).$$

also "sliding window product", non-commutative, kernel "flipped over" (u is transposed in $u+t$). If kernel symmetric: cross-correlation = convolution.

13.3 Convolution via Matrices

Let \mathbf{R} be the signal, the kernel and the output as $m \times n$.

Copy the kernel as columns in the matrix (effecting it by one more time (gives a band matrix (specific of Toeplitz matrix)).

Then the convolution is just a matrix-vector product.

13.3.1 Why use Convolutions in DL

Transforms are usually linear transform + nonlinearity (given in convolution).

Many signals obey translation invariance, so we'd like to have translation invariant features. If the relationship of translation invariant gives us the input-output relation then this is perfect

13.3.2 Border Padding

There are different options to do this

D. (Padding of f)

Mean we extend the image (or each dimension) by p on both sides ($(p+2)p$) and just fill in a constant θ (e.g. 0).

D. (Same Padding)

our definition: padding with zeros = same padding = "same" constant, i.e., 0, and we'll get a tensor of the same dimensions.

D. (Valid Padding)

only retain values from windows that are fully-contained within the support of the signal f (see 2D example below) = valid padding

13.3.3 Backpropagation for Convolutions

Exploits structural sparseness.

D. (Receptive Field Z_i of x_i)

ImageNet: similar to LeNet5, just deeper and using GPU (performance breakthrough)

13.3.4 Inception Module

Now we want to deal with this ever deeper and deeper channels were that the filters at every layer were getting longer and longer and lots of their coefficients were becoming zero (so no information flowing through). So, Arora et al. came up with the idea of an inception module.

Com. One can extend the definition of the receptive field over several layers. The further we go back in layer, the bigger the receptive field becomes due to the nested convolutions. The receptive field may be even the entire image after a few layers. Hence, the convolutions are that is given to the inception module.

13.3.5 Google Inception Network

ImageNet: similar to LeNet5, just deeper and using GPU (performance breakthrough)

13.3.6 Convex Sets

A set $S \subseteq \mathbb{R}^d$ is called convex if

$$x, y \in S, \forall \lambda \in [0, 1]: \quad \lambda x + (1 - \lambda)y \in S.$$

Com. Any point on the line between two points is within the set.

D. (Convex Function)

A function $f: S \rightarrow \mathbb{R}$ defined on a convex set S is called convex if

$$x, y \in S, \lambda \in [0, 1]: \quad f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

Com. Another way to formulate that f is convex function is to say that the graph of f is a convex set.

Com. Every local optimum of a convex function is a global optimum.

D. (Convex Function)

\rightarrow is convex if and only if f is convex

- nonnegative weighted sum

- composition with non-decreasing function, e.g. $c \circ f$

- 1D convex for dimension reduction before convolving with large kernel

- then all these informations are passed to the next layer.

- gradient shortcuts: connect softmax layer at intermediate stage to the output layer (softmax is non-decreasing)

- decomposition of convolution kernels for computational performance

- all the dimensionality reductions improves the efficiency.

D. (Convex Function)

A differentiable function f is called *strictly convex* if

$$x, y \in S, x \neq y: \nabla^2 f(x) - \nabla^2 f(y) > 0 \quad \text{or} \quad \|x - y\|_2^2$$

where $\|\cdot\|$ is any norm. An equivalent condition is the following:

$$x, y: \nabla f(x) - \nabla f(y) = (f(x) - f(y)) \frac{\mu}{\|x - y\|_2^2}$$

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

Com. We can take the insights that we got by the paper of Salakhutdinov et al.

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

D. (Convex Function)

Now we want to show that SGD is $\mathcal{O}(n)$ for a strongly convex function.

Com. The concept of strong convexity extends and parametrizes the notion of local convexity. A strongly convex function is also called *strictly convex*. Note that a convex function is strictly convex, but not strongly convex. So the output nodes do not connect to nodes just to a set of input nodes that are considered "near" (local optimum).

14.11.1 Norm-Based Regularization

$\mathcal{R}(\theta, S) = \mathcal{R}(S) + \Omega(\theta)$,
where Ω is a functional (function from a vector-space to the field over which it's defined) that does not depend on the training data.

D. (L2) Frobenius-Norm Penalty (Weight Decay)

$\Omega(\theta) = \frac{1}{2} \sum_{i,j} \lambda^2 \|W_{i,j}\|^2, \quad \lambda^2 \geq 0$
Com.: It's common practice to only penalize the weights, and not the bias.

The assumption here is that the weights have to be small. So we only allow a big increase in the weights, if it comes at a much bigger increase in performance. Regularization based on the L2-norm is also called weight-decay, as

$$\frac{\partial \mathcal{R}}{\partial w_{i,j}} = \lambda w_{i,j},$$

which means that the weights in the i -th layer get pulled towards zero with "gain" λ^2 . What happens in the gradient-update step is

$$\theta(t+1) = \theta(t) - \eta \nabla_{\theta} \mathcal{R}(\theta, S)$$

$$= (1 - \eta \lambda) \theta(t) - \eta \nabla_{\theta} \mathcal{R}.$$

weight decay step data dep.

and also note that we require $\lambda^2 < 1$.

Let's analyze the weight decay: The Quadratic (Taylor) approximation of \mathcal{R} around the optimal θ^* would be

$$\mathcal{R}(\theta) \approx \mathcal{R}(\theta^*) + \nabla_{\theta} \mathcal{R}(\theta^*)^T (\theta - \theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H} (\theta - \theta^*)$$

$$= \mathcal{R}(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H} (\theta - \theta^*) \quad (*)$$

where $\mathbf{H}_{i,j}$ is the hessian of \mathcal{R} , so

$$(\mathbf{H}_{\theta})_{i,j} = \frac{\partial^2 \mathcal{R}}{\partial \theta_i \partial \theta_j}$$

and \mathbf{H} is the evaluation of \mathbf{H} at θ^* :

$$\mathbf{H} := \mathbf{H}_{\theta^*}(\theta^*).$$

So now we have the upper quadratic approximation of the cost function $(*)$ (so we're assuming it is a parabola and that we know θ^*). Now, let's compute the gradient of that upper approximation of \mathcal{R} in $(*)$:

$$\nabla_{\theta} \mathcal{R}(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H} (\theta - \theta^*) = -\mathbf{H} \theta^* + \mathbf{H} \theta \quad (*)$$

Further, recall that

$$\nabla_{\theta} \Omega = \lambda \theta = \text{diag}(\lambda) \theta$$

So, now, let's set $\nabla_{\theta} \mathcal{R}_{\Omega}$ (with $\nabla_{\theta} \mathcal{R}$ approximated as in $(*)$) equal to zero.

$$\begin{aligned} \nabla_{\theta} \mathcal{R}_{\Omega} &\stackrel{!}{=} 0 \\ \iff -\mathbf{H} \theta^* + \mathbf{H} \theta + \nabla_{\theta} \mathcal{R}(\theta^*) &\stackrel{!}{=} 0 \\ \iff (\mathbf{H} + \text{diag}(\lambda)) \theta &= \mathbf{H} \theta^* \end{aligned}$$

Since both \mathbf{H} and $\text{diag}(\lambda)$ are s.p.d., we can invert their sum

$$\theta = (\mathbf{H} + \text{diag}(\lambda))^{-1} \mathbf{H} \theta^*$$

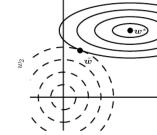
Now, what can we directly see here is that we use no L2-regularization $\theta = \theta^*$. Now, since \mathbf{H} is s.p.d., we can diagonalize it to $\mathbf{H} = \mathbf{Q} \mathbf{A} \mathbf{Q}^T$ where $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_d)$ and plug this in which gives us

$$\begin{aligned} \theta &= (\mathbf{Q} \mathbf{A} \mathbf{Q}^T + \text{diag}(\lambda))^{-1} \mathbf{Q} \mathbf{A} \mathbf{Q}^T \theta^* \\ &= \mathbf{Q} \left(\mathbf{A} + \text{diag}(\lambda) \right)^{-1} \mathbf{A} \mathbf{Q}^T \theta^*. \\ &= \text{diag} \left(\frac{\lambda_1}{\lambda_1 + \lambda_2}, \dots, \frac{\lambda_d}{\lambda_1 + \lambda_d} \right) \theta^*. \end{aligned}$$

So this gives us an idea what happens with θ^* in the directions of the eigenvalues of \mathbf{H} .

If $\lambda > \lambda_i$ effect vanishes along directions in parameter space with large eigenvalues λ_i the weights are almost not reduced. If $\lambda < \lambda_i$ shrinking effect: along the directions in parameter space with small eigenvalues λ_i the weights are shrunk to nearly zero magnitude.

The following picture illustrates this better:



The isometric balls illustrate the regularization loss (L2) for any direction θ for the case that the weights are initialized for the risk of a parabolic risk. $\mathcal{R}(\theta)$ is the point with the least loss for its specific regularization risk. $\mathcal{R}(\theta)$ is the point with the least loss for its specific regularization risk.

So we're essentially clipping the weights.

Actually, for each λ in L2-regularization there is a radius θ that would make the two problems equivalent (if the loss is convex).

Hinton made some nice remarks about this in his paper:

• the gradients do not affect the initial learning (so we'll clip the weights to small at the beginning), so we won't clip the weights.

• down-wards the risk has a large eigenvalue, as the risk increases rapidly. And as we've stated above, the value of w along that dimension is shrunk to zero.

• from right to left (starting at w^*) the risk has a very low eigen-value, and hence w is reduced much more along that dimension.

D. (L1)-Regularization (sparsity-inducing)

$\Omega(\theta) = \sum_{i,j} \lambda^1 |W_{i,j}| = \sum_{i,j} \lambda^1 \sum_{i,j} |w_{i,j}|, \quad \lambda^1 \geq 0$

14.11.1.1 Regularization via Constrained Optimization

An alternative view on regularization is for a given $r > 0$, solve

$$\min_{\theta} \mathcal{R}(\theta).$$

So we're also constraining the size of the coefficients indirectly, by constraining θ to some ball.

The simple optimization approach to this is: projected gradient descent:

$$\theta(t+1) = \theta(t) - \eta \nabla_{\theta} \mathcal{R}.$$

$$\Pi_r(\mathbf{v}) := \min \left\{ 1, \frac{r}{\|\mathbf{v}\|} \right\} \mathbf{v}$$

So we're essentially clipping the weights.

Actually, for each λ in L1-regularization there is a radius θ that would make the two problems equivalent (if the loss is convex).

Hinton made some nice remarks about this in his paper:

• the gradients do not affect the initial learning (so we'll clip the weights to small at the beginning), so we won't clip the weights.

• alternatively, we may just constrain the norm of the incoming weights for each unit (so use row-norms for the weight matrices).

This is a common trick in stabilizing the optimization.

14.11.2 Early Stopping

Gradient descent usually evolves solutions from: simple + robust + complex + sensitive. Hence, it makes sense to stop training early (as soon as validation loss flattens/increases). Also: computationally attractive to stop training early.

Since the weights are initialized to small values (and grow and grow to fit/overfit) we will be kind of clipping/constraining the weight sizes by stopping the learning early.

Advantage of RNNs: better the time component, less memory of past in hidden state.

Given that we're dealing with a dynamical system. We want to clip the hidden activities \mathbf{h}^t with the state of a dynamical system, the discrete time evolution of the hidden state sequence is expressed as a HMM with a non-linearity:

So (as seen previously) we have that

$$\theta(t+1) = \theta(t) - \eta \nabla_{\theta} \mathcal{R}_{\Omega} \approx \theta(t) - \eta \mathbf{H} (\theta(t) - \theta^*).$$

Now, subtracting θ^* on both sides gives

$$\theta(t+1) - \theta^* \approx (\mathbf{I} - \eta \mathbf{H}) (\theta(t) - \theta^*).$$

Now we'll use the same trick as before so that we can diagonalize the distribution of θ it's p.s.d., so $\mathbf{H} = \mathbf{Q} \mathbf{A} \mathbf{Q}^T$. Inserting this gives us:

$$\theta(t+1) - \theta^* \approx (\mathbf{I} - \eta \mathbf{Q} \mathbf{A} \mathbf{Q}^T) (\theta(t) - \theta^*).$$

Now we're making the assumption here is that the weights have to be small. So we only allow a big increase in the weights, if it comes at a much bigger increase in performance. Regularization based on the L2-norm is also called weight-decay, as

$$\frac{\partial \mathcal{R}}{\partial w_{i,j}} = \lambda w_{i,j},$$

which means that the weights in the i -th layer get pulled towards zero with "gain" λ^2 . What happens in the gradient-update step is

$$\begin{aligned} \theta(t+1) &= \theta(t) - \eta \nabla_{\theta} \mathcal{R}_{\Omega} \approx \theta(t) - \eta \mathbf{H} (\theta(t) - \theta^*) \\ &= (1 - \eta \lambda) \theta(t) - \eta \nabla_{\theta} \mathcal{R}. \end{aligned}$$

weight decay step data dep.

and also note that we require $\lambda^2 < 1$.

Let's analyze the weight decay: The Quadratic (Taylor) approximation of \mathcal{R} around the optimal θ^* would be

$$\mathcal{R}(\theta) \approx \mathcal{R}(\theta^*) + \nabla_{\theta} \mathcal{R}(\theta^*)^T (\theta - \theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H} (\theta - \theta^*)$$

$$= \mathcal{R}(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H} (\theta - \theta^*) \quad (*)$$

where \mathbf{H} is the hessian of \mathcal{R} , so

$$(\mathbf{H}_{\theta})_{i,j} = \frac{\partial^2 \mathcal{R}}{\partial \theta_i \partial \theta_j}$$

and \mathbf{H} is the evaluation of \mathbf{H} at θ^* :

$$\mathbf{H} := \mathbf{H}_{\theta^*}(\theta^*).$$

So now we have the upper quadratic approximation of the cost function $(*)$ (so we're assuming it is a parabola and that we know θ^*). Now, let's compute the gradient of that upper approximation of \mathcal{R} in $(*)$:

$$\nabla_{\theta} \mathcal{R}(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T \mathbf{H} (\theta - \theta^*) = -\mathbf{H} \theta^* + \mathbf{H} \theta \quad (*)$$

Further, recall that

$$\nabla_{\theta} \Omega = \lambda \theta = \text{diag}(\lambda) \theta$$

So, now, let's set $\nabla_{\theta} \mathcal{R}_{\Omega}$ (with $\nabla_{\theta} \mathcal{R}$ approximated as in $(*)$) equal to zero.

$$\begin{aligned} \nabla_{\theta} \mathcal{R}_{\Omega} &\stackrel{!}{=} 0 \\ \iff -\mathbf{H} \theta^* + \mathbf{H} \theta + \nabla_{\theta} \mathcal{R}(\theta^*) &\stackrel{!}{=} 0 \\ \iff (\mathbf{H} + \text{diag}(\lambda)) \theta &= \mathbf{H} \theta^* \end{aligned}$$

Since both \mathbf{H} and $\text{diag}(\lambda)$ are s.p.d., we can invert their sum

$$\theta = (\mathbf{H} + \text{diag}(\lambda))^{-1} \mathbf{H} \theta^*$$

Now, what can we directly see here is that we use no L2-regularization $\theta = \theta^*$. Now, since \mathbf{H} is s.p.d., we can diagonalize it to $\mathbf{H} = \mathbf{Q} \mathbf{A} \mathbf{Q}^T$ where $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_d)$ and plug this in which gives us

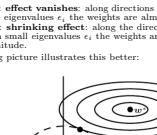
$$\begin{aligned} \theta &= (\mathbf{Q} \mathbf{A} \mathbf{Q}^T + \text{diag}(\lambda))^{\frac{1}{2}} \mathbf{Q} \mathbf{A} \mathbf{Q}^T \theta^* \\ &= \mathbf{Q} \left(\mathbf{A} + \text{diag}(\lambda) \right)^{\frac{1}{2}} \mathbf{A} \mathbf{Q}^T \theta^*. \\ &= \text{diag} \left(\frac{\lambda_1}{\lambda_1 + \lambda_2}, \dots, \frac{\lambda_d}{\lambda_1 + \lambda_d} \right) \theta^*. \end{aligned}$$

So this gives us an idea what happens with θ^* in the directions of the eigenvalues of \mathbf{H} .

If $\lambda > \lambda_i$ effect vanishes along directions in parameter space with large eigenvalues λ_i the weights are almost not reduced.

If $\lambda < \lambda_i$ shrinking effect: along the directions in parameter space with small eigenvalues λ_i the weights are shrunk to nearly zero magnitude.

The following picture illustrates this better:



The isometric balls illustrate the regularization loss (L2) for any direction θ for the case that the weights are initialized for a parabolic risk. $\mathcal{R}(\theta)$ is the point with the least loss for its specific regularization risk.

So we're essentially clipping the weights.

Actually, for each λ in L2-regularization there is a radius θ that would make the two problems equivalent (if the loss is convex).

Hinton made some nice remarks about this in his paper:

• the gradients do not affect the initial learning (so we'll clip the weights to small at the beginning), so we won't clip the weights.

• alternatively, we may just constrain the norm of the incoming weights for each unit (so use row-norms for the weight matrices).

This is a common trick in stabilizing the optimization.

14.11.2 Early Stopping

Gradient descent usually evolves solutions from: simple + robust + complex + sensitive. Hence, it makes sense to stop training early (as soon as validation loss flattens/increases). Also: computationally attractive to stop training early.

Since the weights are initialized to small values (and grow and grow to fit/overfit) we will be kind of clipping/constraining the weight sizes by stopping the learning early.

Advantage of RNNs: better the time component, less memory of past in hidden state.

Given that we're dealing with a dynamical system. We want to clip the hidden activities \mathbf{h}^t with the state of a dynamical system, the discrete time evolution of the hidden state sequence is expressed as a HMM with a non-linearity:

$$\theta = (\mathbf{U}, \mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c})$$

$$\mathbf{h}^t = \mathbf{F}(\mathbf{h}^{t-1}, \mathbf{x}^t; \theta), \quad \mathbf{h}^0 = \mathbf{o}$$

$$F = \sigma, \text{ tanh}, \dots, \text{ logistic}, \text{ ReLU}, \dots$$

$$\mathbf{F}(\mathbf{h}, \mathbf{x}) = \mathbf{W} \mathbf{h} + \mathbf{U} \mathbf{x} + \mathbf{b}$$

$$\mathbf{y}^t = \mathbf{H}(\mathbf{h}^t) = \mathbf{V}(\mathbf{h}^t) + \mathbf{c}$$

This is just because the Jacobian of the gradient map is the Hessian \mathbf{H} from before.

This is just because the Jacobian of the gradient map is the Hessian \mathbf{H} from before.

So (as seen previously) we have that

$$\theta(t+1) = \theta(t) - \eta \nabla_{\theta} \mathcal{R}_{\Omega} + \eta \nabla_{\theta} \mathcal{R}_{\Omega}(\theta = \mathbf{H}(\theta = \theta^*))$$

This ensures that the gradient norm is never greater than γ_{max} .

However, when having vanishing gradients over time, then this is not a good idea.

• the gradient norm would usually be 0, we'd never learn to predict the next hidden state.

• the gradient norm would be very large, it's a huge simplification that works.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

• the gradient norm would be 0, we'd have to learn to predict the next hidden state.

17.2.1 Linear Autoencoder

D. (Linear Autoencoder)
A linear autoencoder just consists of two linear maps: an encoder $\mathbf{C} \in \mathbb{R}^{m \times d}$ and a decoder $\mathbf{D}^T \in \mathbb{R}^{d \times m}$. The objective it minimizes is then:

$$R(\theta) = \frac{1}{2} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\mathbf{C}\mathbf{x}_i\|_2^2.$$

So it's a NN with one hidden layer (no biases and linear activation functions) which will contain the compressed representation $\mathbf{z} = \mathbf{C}\mathbf{x} + \mathbf{b}$.

D. (Linear Autoencoder with Coupled Weights)

Then, we define $\mathbf{D} := \mathbf{C}^T$.

D. (Singular Value Decomposition)

Recall that the SVD of a data matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n \end{bmatrix}$$

is of the following form:

$$\mathbf{X} = \mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_{\min(n, k)}) \mathbf{V}^T$$

And the matrices \mathbf{U} and \mathbf{V} are orthogonal - so we have an orthogonal basis. Further recall that via the SVD we can get the best rank m approximation of a linear mapping. It also is a decomposition that preserves the energy (or the energy) of the data for a predefined number of desired basis vectors to represent it.

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}_m^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

$\mathbf{D}\mathbf{C} = \mathbf{U}_m \mathbf{U}_m^T \mathbf{U}_m \mathbf{V}_m^T = \mathbf{U}_m [\mathbf{I}_m \quad \mathbf{0}] \mathbf{V}_m^T = \mathbf{U}_m [\mathbf{I}_m \quad \mathbf{0}] \mathbf{V}_m^T$

And as we know from the Eckhart-Young theorem $\mathbf{X} = \mathbf{U}_m \mathbf{V}_m^T$ is the best m -dimensional approximation of the original data \mathbf{X} .

Now, since $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ that means that we can do weight sharing between the decoder and encoder network, since $\mathbf{C} = \mathbf{D}^T$.

Another thing to note is that the solution is not unique! For any invertible matrix $\mathbf{A} \in GL(m)$ we have:

$$\frac{(\mathbf{U}_m \mathbf{A})^T (\mathbf{A}^T)^T}{\mathbf{D}} = \mathbf{U}_m \mathbf{A}^T$$

Now, restricting through weight sharing that $\mathbf{D} = \mathbf{C}^T$ will enforce that

$$\mathbf{A}^T = \mathbf{A}^T$$

because $\mathbf{A} \in GL(m)$ (orthogonal group, rotation matrices). Then the mapping $\mathbf{x} \rightarrow \mathbf{z}$ is deterministic (up some rotation that we do in-between, rotation and its inverse).

D. (Principal Component Analysis)

This is a very simple problem through PCA. First, we center the data (pre-processing) as follows:

$$\mathbf{x}_i \rightarrow \mathbf{x}_i - \mathbf{x}_i^{\text{mean}}$$

Then we define

$$\mathbf{S} = \mathbf{XX}^T$$

which is the sample covariance matrix. And then, in order to get \mathbf{U} we just do the singular value decomposition of \mathbf{S} . If we relate to the SVD of \mathbf{X} we can see that

$$\mathbf{S} = \mathbf{U}\mathbf{\Sigma}\mathbf{U}^T = \mathbf{U}^2\mathbf{\Sigma}^2\mathbf{U}^T.$$

So, the columns of \mathbf{U} are the eigenvectors of the covariance matrix. And \mathbf{U}_m^T is the orthogonal projection onto its principal components of \mathbf{S} .

Note that if we wanted to get the \mathbf{V} we'd just do the PCA with $\mathbf{S} = \mathbf{X}^T\mathbf{X}$.

17.2.2 Non-Linear Autoencoders
Non-linear autoencoders allow us to learn powerful non-linear generalization of the PCA.

D. (Non-Linear Autoencoder) contains many hidden layers with nonlinearities and want to learn “overcomplete” codes.

D. (Code Sparseness) e.g. $\|\mathbf{z}\|_1 = \lambda \|\mathbf{z}\|_2$.

D. (Contractive Autoencoders) $R(\mathbf{z}) = \lambda \|\mathbf{z}\|_2^2$. This penalizes the Jacobian and generalizes weight decay (cf. Rifai et al., 2011).

D. (Denoising Autoencoders)

Autoencoders allow us to separate the signal from noise: Denoising autoencoders are the first step of the original data representation that are robust under noise.

D. (Denoising Autoencoder)

we perturb the inputs

$$\mathbf{x} \rightarrow \mathbf{x}_n$$

where η is a random noise vector, e.g. additive (white) noise

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

and instead of the original objective, we minimize the following

$$E_{\mathbf{p}(\mathbf{x})}[\mathbf{I}(\mathbf{x}, (\mathbf{H}(\mathbf{x}) \mathbf{p}(\mathbf{x})))]$$

The hope is that we'll achieve *de-noising*, which happens if

$$\|\mathbf{x} - \mathbf{H}(\mathbf{x})\mathbf{p}(\mathbf{x})\|^2 < \|\mathbf{x} - \mathbf{x}_n\|^2$$

So this would mean that the reconstruction error of the noisy data is less than the error we created by the noise we added (then the denoising works).

D. (Factor Analysis)

17.3.1 Latent Variable Analysis

Latent Variable Analysis provides a generic way of defining probabilities for a data set given the so-called latent variable models. They usually work as follows:

1. Define a *latent variable* \mathbf{z} , with a distribution $\mathbf{p}(\mathbf{z})$.

2. Define *conditional models* for the observable \mathbf{x} conditioned on the latent variable \mathbf{z} .

3. Construct the *observed model* by integrating/summing out the latent variables

$$p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z}) \mathbf{p}(\mathbf{z}) d\mathbf{z} = \left[\int p(\mathbf{x} | \mathbf{z}) d\mathbf{z} \right] \mathbf{p}(\mathbf{z}) = \text{constant}$$

and the empirical co-variance matrix is

$$\mathbf{S} = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{k} \mathbf{XX}^T$$

Then, the log-likelihood of the data \mathbf{X} , given \mathbf{A} can be written as:

$$\log(P(\mathbf{X} | \mathbf{A})) = -\frac{k}{2} (\mathbf{Tr}(\mathbf{SA}^{-1}) - \log(\det(\mathbf{A}))) + \text{const.}_{\text{I.T. of } \mathbf{A}}$$

The idea of latent variable models is very similar to the one of autoencoders. The idea is to have some

$$\cdot \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

• and want to embed it into \mathbb{R}^k ($k \ll d$)

• so we'll use $\mathbf{x} \in \mathbb{R}^k$ (latent-space)

• and look at the conditional probabilities $\mathbf{p}(\mathbf{x} | \mathbf{z})$ for some \mathbf{z} computing the matrices \mathbf{B} and \mathbf{C}

Ex. (Gaussian Mixture Models GMMS)

$\mathbf{x} \in \{1, \dots, K\}$, $\mathbf{p}(\mathbf{x})$ = mixing proportions

$p(\mathbf{x} | \mathbf{z})$: conditional densities (Gaussians for GMMS)

The idea of latent variable models is very similar to the one of autoencoders. The idea is to have some

$$\cdot \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

• and want to embed it into \mathbb{R}^k ($k \ll d$)

• so we'll use $\mathbf{x} \in \mathbb{R}^k$ (latent-space)

• and look at the conditional probabilities $\mathbf{p}(\mathbf{x} | \mathbf{z})$ for some \mathbf{z}

Now, let's compute the matrix gradients w.r.t. \mathbf{A} to know the equations that we need to compute the maximum likelihood:

$$\nabla_{\mathbf{A}} \operatorname{Tr}(\mathbf{SA}^{-1}) = -\mathbf{A}^{-1} \mathbf{S}$$

$\nabla_{\mathbf{A}} \log(\det(\mathbf{A})) = \mathbf{A}^{-1}$

Now, setting the gradient of the log-likelihood to zero gives us the following condition:

$$\nabla_{\mathbf{A}} \log(P(\mathbf{X} | \mathbf{A})) = 0 \iff \mathbf{SA}^{-1} = \mathbf{I}.$$

So, the MLE for \mathbf{A} is just $\mathbf{A} = \mathbf{S}$.

But recall that what we want is not \mathbf{A} , but we want \mathbf{W} and \mathbf{E} . The reconstruction is done via a linear map \mathbf{W} and then different gaussian noises are added to the reconstructed vector ($\mathbf{w} + \mathbf{e}$).

So the latent variable $\mathbf{z} \in \mathbb{R}^k$ where

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

and we have a linear *observation model* for $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} = \mathbf{w} + \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Further note that

• \mathbf{w} and \mathbf{e} are $\mathbf{independent}$

• so few factors account for the dependencies between many observables

• The vector \mathbf{w} is computed through MLE on the training set

$$\mathbf{w} = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i$$

is of the following form:

$$\mathbf{w} = \mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_{\min(n, k)}) \mathbf{V}^T$$

And the matrices \mathbf{U} and \mathbf{V} are orthogonal - so we have an orthogonal basis. Further recall that via the SVD we can get the best rank m approximation of a linear mapping. It also is a decomposition that preserves the energy (or the energy) of the data for a predefined number of desired basis vectors to represent it.

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**

D. (Optimal Linear Compression)

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times m}} \|\mathbf{X} - \mathbf{W}\|_F^2 = \mathbf{W} \operatorname{diag}(\sigma_1, \dots, \sigma_m) \mathbf{V}^T$$

where the subscript m refers to the matrices of the SVD pruned to m columns. This means that a linear autoencoder with m hidden units can reconstruct the input \mathbf{X} with $\mathbf{E}(\mathbf{X})$. However, the auto-encoder can achieve the result of the SVD.

T. (Given the SVD of the data \mathbf{X} is $\mathbf{U} \operatorname{diag}(\sigma_1, \dots, \sigma_n) \mathbf{V}^T$.

The choice $\mathbf{C} = \mathbf{U}_m^T$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of two-layer linear auto-encoder with m hidden units. **Proof.**