

1 Probability

Sum Rule $P(X = x_i) = \sum_j p(X = x_i, Y = y_i)$
Product rule $P(X, Y) = P(Y|X)P(X)$
Independence $P(X, Y) = P(X)P(Y)$
Bayes' Rule $P(Y|X) = \frac{P(X|Y)P(Y)}{\sum_{i=1}^n P(X|Y_i)P(Y_i)}$
Cond. Ind. $X \perp\!\!\!\perp Y|Z \Rightarrow P(X, Y|Z) = P(X|Z)P(Y|Z)$
Cond. Ind. $X \perp\!\!\!\perp Y|Z \Rightarrow P(X|Y, Z) = P(X|Z)$
 $E[X] = \int_X t \cdot f_X(t) dt = \mu_X$
 $\text{Cov}(X, Y) = E_{x,y}[X - E_x[X](Y - E_y[Y])]$
 $\text{Cov}(X, Y) := \text{Cov}(X, X) = \text{Var}[X]$
 X, Y independent $\Rightarrow \text{Cov}(X, Y) = 0$
 $\text{XX}^T \geq 0$ (symmetric positive semidefinite)
 $\text{Var}[X] = E[X^2] - E[X]^2$
 $\text{Var}[\mathbf{A}X] = \mathbf{A} \text{Var}[X] \mathbf{A}^T$ $\text{Var}[aX + b] = a^2 \text{Var}[X]$
 $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + 2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$
 $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$
 $\frac{\partial}{\partial t} P(X \leq t) = \frac{\partial}{\partial t} F_X(t) = f_X(t)$ (derivative of c.d.f. is p.d.f.)
 $f_{\alpha|Y}(z) = \frac{1}{\alpha} f_Y(\frac{z}{\alpha})$
T: The moment generating function (MGF) $\psi_X(t) = E[e^{tX}]$ characterizes the distr. of a rv
 $B(p): p^t + (1-p)^{1-t}$ $N(\mu, \sigma): \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$
 $\text{Bin}(n, p): (pe + (1-p))^n$ $\text{Gam}(\alpha, \beta): (\frac{\alpha}{\alpha+\beta})^\alpha$ for $t < 1/\beta$
 $\text{Pois}(\lambda): e^{\lambda(e^{t-1})}$
T: If X_1, \dots, X_n are indep. rvs with MGFs $M_{X_i}(t) = E[e^{tX_i}]$, then the MGF of $Y = \sum_{i=1}^n a_i X_i$ is $M_Y(t) = \prod_{i=1}^n M_{X_i}(a_i t)$.
T: Let X, Y be indep., then the p.d.f. of $Z = X + Y$ is the conv. of the p.d.f. of X and Y : $f_Z(z) = \int_R f_X(x) f_Y(z-x) dt = \int_R f_X(x) f_Y(z-x) dt$
D: $\text{KL}(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(\frac{q(x)}{p(x)}) = \sum_{x \in \mathcal{X}} p(x) \log(\frac{p(x)}{q(x)}) \geq 0$
 $\text{KL}(p, q) = -E_{x \sim p} [\log(\frac{q(x)}{p(x)})] = E_{x \sim p} [\log(\frac{p(x)}{q(x)})] \geq 0$
 $KL(X; p, q) = H(p, q) - H(X)$, where H uses p .
Com. The second formulation clearly shows why p is the minimizer of the cross-entropy (or hence: the maximizer of the likelihood).
Com. Usually, q is the approximation of the unknown p .
D: **(Kullback-Leibler Divergence)**
For discrete probability distributions p and q defined on the same probability space, the KL-divergence between p and q is defined as
 $KL(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(\frac{q(x)}{p(x)}) = \sum_{x \in \mathcal{X}} p(x) \log(\frac{p(x)}{q(x)}) \geq 0$
 $\text{KL}(p, q) = -E_{x \sim p} [\log(\frac{q(x)}{p(x)})] = E_{x \sim p} [\log(\frac{p(x)}{q(x)})] \geq 0$
 $KL(X; p, q) = H(p, q) - H(X)$, where H uses p .
The KL-divergence is defined only if $\forall x: q(x) = 0 \Rightarrow p(x) = 0$ (absolute continuity). Whenever $p(x)$ is zero the contribution of the corresponding term is interpreted as zero because it's 0.
 $\frac{\partial}{\partial t} P(X \leq t) = \frac{\partial}{\partial t} F_X(t) = f_X(t)$ (derivative of c.d.f. is p.d.f.)
 $\exists h \in \mathbb{R}: \max_{x \in R} |f(x) - h(x)| = \|f - h\|_\infty, K < \epsilon$.
T: In ML it is a measure of the amount of information lost, when q (model) is used to approximate p (true).
Com. $KL(p, q) = 0 \iff p \equiv q$.
Com. Note that the KL-divergence is not symmetric!
D: **(Jensen-Shannon Divergence)**
 $R(P, Q) = \frac{1}{2}KL(P, M) + \frac{1}{2}KL(Q, M) \in [0, \log(n)]$ $M = \frac{1}{2}(P + Q)$
T: The JSO is symmetric!
Com. The JSO is a symmetrized and smoothed version of the KL-divergence.
D: **(Smooth Function f in C[∞](R))**
A function class $\mathcal{H} \subseteq C(\mathbb{R}^d)$ is dense in $C(\mathbb{R}^d)$, iff Law of large numbers guarantees: $\mathcal{R}(F, S_n) \xrightarrow{n \rightarrow \infty} \mathcal{R}(F)$ (but this doesn't help us to move from some $F \rightarrow F^*$).
In order to ultimately get from \hat{F} to F^* , which minimizes $\mathcal{R}(F^*)$, we need to restrict F such that $\mathcal{R}(\hat{F}, S_n) \xrightarrow{n \rightarrow \infty} \mathcal{R}(F^*)$. This will prevent overfitting. However it's not so easy to determine how to constrain F such that $\hat{F} \rightarrow F^*$.
D: **(Regularized Empirical Risk)**
 $\mathcal{R}_\theta(F, S_n) := \mathcal{R}(F, S_n) + \lambda \Omega(\|F\|)$
T: **(10.1.3 — Regularization Approaches)**
Weight decay, data augmentation, noise robustness (to weights, inputs, labels (compensates for errors in labeling of data)), semi-supervised learning, dropout, early stopping (prevent overfitting through long training), multi-task learning (to learn general representation or use more data), parameter sharing (CNNs), ensembles (reduces variance, same bias, costly)
Some regularization methods can be proven to be equivalent. However: only in the limit. Thus, it's good to use combinations of them in finite horizons.
T: **(Montufar et al., 2014)** If we process n -dim. inputs through L ReLU layers with widths $m_1, \dots, m_L \in \mathcal{O}(m)$. Then \mathbb{R}^n can be partitioned into at most $R(m, L)$ layers:

$$R(m, L) = \sum_{i=0}^{\min(m,n)} \binom{m}{i}$$

C: If classes $m \leq n$ we have $R(m) = 2^m$ (exponential growth).
C: For any input size n we have $R(m) \in \mathcal{O}(n^m)$ (polynomial slowdown in number of cells, limited by the input space dimension).
T: **(10.5.3 — Deep Combination of ReLU's)**
Question: Process n inputs through L ReLU layers with widths $m_1, \dots, m_L \in \mathcal{O}(m)$. Into how many ($= R(m, L)$) cells can \mathbb{R}^n be maximally partitioned?
T: **(Montufar et al., 2014)** If we process n -dim. inputs through L ReLU layers with widths $m_1, \dots, m_L \in \mathcal{O}(m)$. Then \mathbb{R}^n can be partitioned into at most $R(m, L)$ layers:

$$R(m, L) = \sum_{i=0}^{\min(m,n)} \binom{m}{i}$$

11 Backpropagation
Learning in neural networks is about *gradient-based optimization* (with very few exceptions). So what we'll do is compute the gradient of the objective (empirical risk) with regards to the parameters θ :
 $\nabla_\theta \mathcal{R} = \left(\frac{\partial \mathcal{R}}{\partial \theta_1}, \dots, \frac{\partial \mathcal{R}}{\partial \theta_d} \right)^\top$.
11.1 — Comp. of Gradient via Backpropagation
A good way to compute $\nabla_\theta \mathcal{R}$ is by exploiting the computational structure of the network through the so-called *backpropagation*. The basic steps of *backpropagation* are as follows:
1. Forward-Pass: Perform a forward pass (for a given training input \mathbf{x}), compute all the activations for all units
2. Compute the gradient of R w.r.t. the output layer activations (for a given target \mathbf{y}) (even though we're not interested directly in these gradients, they'll simplify the computations of the gradients in step 4.)
3. Iteratively propagate the activation gradient information from the outputs of the previous layer to the inputs of the current layer.
4. Compute the local gradients of the activations w.r.t. the weights and biases
Since NNs are based on the composition of functions we'll inevitably need the *chain rule*.
T: **(1-D Chain Rule)** $y = f(g(x))$

$$(f \circ g)' = (f' \circ g) \cdot g' \quad \frac{d(f \circ g)}{dx} \Big|_{x=0} = \frac{df}{dx} \Big|_{x=g(0)} \cdot \frac{dg}{dx} \Big|_{x=0}$$

D: **(Jacobi Matrix of a Map)** The Jacobian J_F of a map $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is defined as

$$J_F := \begin{pmatrix} -(\nabla F_1)^\top \\ -(\nabla F_2)^\top \\ \vdots \\ -(\nabla F_m)^\top \end{pmatrix} = \begin{pmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_1}{\partial x_2} & \dots & \frac{\partial F_1}{\partial x_n} \\ \frac{\partial F_2}{\partial x_1} & \frac{\partial F_2}{\partial x_2} & \dots & \frac{\partial F_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_m}{\partial x_1} & \frac{\partial F_m}{\partial x_2} & \dots & \frac{\partial F_m}{\partial x_n} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

So each component function $F_i: \mathbb{R}^n \rightarrow \mathbb{R}$ of F , for $i \in \{1, \dots, m\}$ has a constant partial derivative. Putting these gradients together as rows of a matrix gives us the *Jacobian* of F .
So we have that

$$(J_F)_{ij} := \frac{\partial F_j}{\partial x_i}$$

T: **(Jacobi Matrix Chain Rule)**

$$G: \mathbb{R}^n \xrightarrow{G} \mathbb{R}^q \quad H: \mathbb{R}^q \xrightarrow{H} \mathbb{R}^m$$

$$F := H \circ G$$

$$F: \mathbb{R}^n \xrightarrow{G} \mathbb{R}^q \xrightarrow{H} \mathbb{R}^m$$

$$\mathbf{x} \xrightarrow{G} \mathbf{z} := G(\mathbf{x}) \xrightarrow{H} \mathbf{y} := H(\mathbf{z}) = H(G(\mathbf{x}))$$

Chain-rule for a single component-function $F_i: \mathbb{R}^n \rightarrow \mathbb{R}$ of F :

$$\frac{\partial F_i}{\partial x_j} \Big|_{\mathbf{x}=\mathbf{x}_0} = \frac{\partial (H \circ G)_i}{\partial x_j} \Big|_{\mathbf{x}=\mathbf{x}_0} = \sum_{k=1}^q \frac{\partial H_i}{\partial z_k} \Big|_{\mathbf{z}=G(\mathbf{x}_0)} \cdot \frac{\partial G_k}{\partial x_j} \Big|_{\mathbf{x}=\mathbf{x}_0}$$

This gives us the following lemma for Jacobi matrices of compositions of functions
L: **(Jacobi Matrix Chain Rule)**

$$J_F|_{\mathbf{x}=\mathbf{x}_0} = J_{H \circ G}|_{\mathbf{x}=\mathbf{x}_0} = J_H|_{\mathbf{z}=G(\mathbf{x}_0)} \cdot J_G|_{\mathbf{x}=\mathbf{x}_0}$$

T: **(Goodfellow, 2013)** Maxout networks with two maxout units are universal function approximators.
T: **(Rectification Networks)**

$$f(x) \approx f(x) + (\nabla f(x))^{-1} \nabla f(x)$$

D: **(Rectified Linear Unit (ReLU))**

$$f(x) := \max(0, x) = \text{ReLU}(x) \in [0, \infty]$$

$$f'_+(x) := \mathbb{1}_{\{x > 0\}}$$

$$\nabla_x f(\mathbf{x}) + \mathbf{J}_{+}(\mathbf{x}) := \text{diag}(\mathbb{1}_{\{x > 0\}})$$

9 Approximation Theory
9.1 — Compositional Models
We want to learn: $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$.
Learning: Now we reduce this task to learning a function F in some parameter space \mathcal{R}^d that approximates F^* well.
 $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathcal{F} := \{F \in \mathcal{R}^d : \theta \in \mathcal{R}^d\}$
DL: the composition of simple functions can give rise to very complex functions.
 $F: \mathbb{R}^n \xrightarrow{G_1} \mathbb{R}^m \xrightarrow{G_2} \mathbb{R}^n \xrightarrow{G_3} \mathbb{R}^m$
 $F(x, \theta) = G_1(\dots G_2(G_1(\mathbf{x}; \theta_1); \theta_2); \dots; \theta_d)$
9.2 — Compositions of Maps
D: **(Linear Function)** A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear function if the following properties hold

- $\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^n: f(\mathbf{x} + \mathbf{x}') = f(\mathbf{x}) + f(\mathbf{x}')$
- $\forall \mathbf{x} \in \mathbb{R}^n: f(a\mathbf{x}) = af(\mathbf{x})$

T: **(4.3 — Vector-by-Vector)**
A, C, D, a, b, e not a function of \mathbf{x} ,
 $f = \mathbf{f}(\mathbf{x}), \mathbf{g} = \mathbf{g}(\mathbf{x}), \mathbf{h} = \mathbf{h}(\mathbf{x}), u = u(\mathbf{x}), v = v(\mathbf{x})$
 $\frac{\partial}{\partial \mathbf{x}} [u(\mathbf{x})v(\mathbf{x})] = u(\mathbf{x}) \frac{\partial}{\partial \mathbf{x}} v(\mathbf{x}) + v(\mathbf{x}) \frac{\partial}{\partial \mathbf{x}} u(\mathbf{x})$
 $\frac{\partial}{\partial \mathbf{x}} [u(\mathbf{x})v(\mathbf{x})] = \frac{\partial u(\mathbf{x})}{\partial \mathbf{x}} \mathbf{g}(\mathbf{x}) + \frac{\partial v(\mathbf{x})}{\partial \mathbf{x}} \mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) \mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x}) \mathbf{f}(\mathbf{x})$
 $\frac{\partial}{\partial \mathbf{x}} [(\mathbf{A} + \mathbf{B})^\top \mathbf{C}(\mathbf{D} + \mathbf{E})] = \mathbf{D}^\top \mathbf{C}^\top (\mathbf{A} + \mathbf{B}) + \mathbf{A}^\top \mathbf{C}^\top (\mathbf{D} + \mathbf{E})$
 $\frac{\partial}{\partial \mathbf{x}} [(\mathbf{A} + \mathbf{B})^\top \mathbf{C}(\mathbf{D} + \mathbf{E})] = \mathbf{D}^\top \mathbf{C}^\top (\mathbf{A} + \mathbf{B}) + \mathbf{A}^\top \mathbf{C}^\top (\mathbf{D} + \mathbf{E})$
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{A}^\top \mathbf{B}] = \mathbf{A}^\top \frac{\partial \mathbf{B}}{\partial \mathbf{x}}$
4.4 — Scalar-by-Matrix
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X} \mathbf{b}] = \mathbf{a}^\top$
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{b}^\top$
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{b}^\top$
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{a}^\top$
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{a}^\top$
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{a}^\top$
4.5 — Vector-by-Matrix (Generalized Gradient)
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{X} \mathbf{a}] = \mathbf{X}^\top$
5 Information Theory
D: **(Entropy)** Let X be a random variable distributed according to $P(X)$. Then the entropy of X is
 $H(X) = \mathbb{E}[-\log(P(X))] = -\sum_{x \in \mathcal{X}} p(x) \log(p(x)) \geq 0$
It describes the expected information content $I(X)$ of X .
D: **(Cross-Entropy)** For distributions p and q over a given set is $H(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(q(x)) = \mathbb{E}_{x \sim p} [-\log(q(x))] \geq 0$.
D: **(Universe of Ridge Functions for some $\sigma: \mathbb{R} \rightarrow \mathbb{R}$)**
 $\mathcal{G}_\sigma := \{g \mid g(x) = \sigma(w^\top \mathbf{x} + b), \mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}, \sigma: \mathbb{R} \rightarrow \mathbb{R}\}$
D: **(Universes of Continuous Ridge Functions)**
 $\mathcal{G}^n := \cup_{\sigma \in C(\mathbb{R})} \mathcal{G}_\sigma$
Com. The second formulation clearly shows why $g = p$ is the minimizer of the cross-entropy (or hence: the maximizer of the likelihood).
Com. g is the approximation of the unknown p .
D: **(Kullback-Leibler Divergence)**
For discrete probability distributions p and q defined on the same probability space, the KL-divergence between p and q is defined as
 $KL(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(\frac{q(x)}{p(x)}) = \sum_{x \in \mathcal{X}} p(x) \log(\frac{p(x)}{q(x)}) \geq 0$
 $KL(p, q) = -E_{x \sim p} [\log(\frac{q(x)}{p(x)})] = E_{x \sim p} [\log(\frac{p(x)}{q(x)})] \geq 0$
 $KL(X; p, q) = H(p, q) - H(X)$, where H uses p .
Com. The second formulation clearly shows why $g = p$ is the minimizer of the cross-entropy (or hence: the maximizer of the likelihood).
Com. g is the approximation of the unknown p .
D: **(Piecewise Linear Functions and Half Spaces)**
So the ReLU and the AVU define a *piecewise linear function* with 2 pieces. Hereby, \mathbb{R}^n is partitioned into two open half spaces (and a border face):

- $H^+ := \{ \mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b > 0 \} \subset \mathbb{R}^n$
- $H^- := \{ \mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b < 0 \} \subset \mathbb{R}^n$
- $H^0 := \{ \mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b = 0 \} = \mathbb{R}^n - H^+ - H^- \subset \mathbb{R}^n$

Further note that

- $g_{-}(+)(H^0) = g_{-1}(H^0) = 0$
- $g_{-1,+}(H^-) = 0$
- $g_{-1,-}(H^+) = 0$

So the *training risk* is the *expected risk* under the empirical distribution induced by the sample S_N .

12.6 — Backpropagation for Convolutions

Exploits structural sparseness.

D. (Receptive Field T_i^l of x_i^l)

The receptive field T_i^l of node x_i^l is defined as $T_i^l := \{j \mid W_{ij}^l \neq 0\}$ where \mathbf{W}^l is the Toeplitz matrix of the convolution at layer l .

Com. Hence, the receptive field of a node x_i^l are just nodes which are connected to it and have a non-zero weight.

Com. One may 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 have to be small.

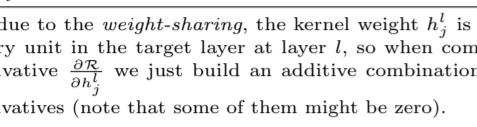
This is what gave rise to the inception module:

D. (Dimension Reduction) m channels of a $1 \times 1 \times k$ convolution $m \leq k$:

$$x^{l-1}ij = \sigma(\mathbf{W}x_{ij}), \quad \mathbf{W} \in \mathbb{R}^{m \times k}.$$

So it uses a 1×1 filter over the k input channels (which is actually not a convolution), aka "network within a network".

12.13.5 — Google Inception Network



The Google Inception Network uses many layers of this inception module along with some other tricks

- dimensionality reduction through the inception modules
- convolutions at various sizes, as different filter sizes turned out to be useful
- a max-pooling of the previous layer, and a dimensionality reduction of the result
- composition for dimension reduction before convolving with large kernels
- then all these informations are passed to the next layer
- gradient shortcuts: connect softmax layer at intermediate stages to have the gradient flow until the beginnings of the network
- decomposition of convolution kernels for computational performance
- in all the dimensionality reductions improved the efficiency

12.14 — Networks Similar to CNNs

D. (Locally Connected Network) A locally connected network has the same connections that a CNN would have, however, the parameters are not shared. So the output nodes do not connect to all nodes, just to a set of input nodes that are considered "near" (locally connected).

12.15 — Comparison of #Parameters (CNNs, FC, LC)

Ex. Input image $m \times n \times c$ (c number of channels)

K convolution kernels: $p \times q$ (valid padding and stride 1)

output dimensions: $(m-p+1) \times (n-q+1) \times K$

#parameters: CNN: $K(pqc + 1)$

#parameters of fully-conn. NN with same number of outputs as CNN: $mnc(m-p+1)(n-q+1)+1$

#parameters of locally-conn. NN with same connections as CNN: $pqc(m-p+1)(n-q+1)+K$

Ex. Assume we have an $m \times n$ image (with one channel).

And we convolve it with a filter $(2p+1) \times (2q+1)$

Then the convolved image has dimensions (assuming stride 1)

- valid padding (only where it's defined): $(m-2p) \times (n-2q)$
- padded (extend image with constant): $m \times n$ where the extended image has size $(m+2p) \times (n+2q)$.

12.8 — Typical Convolutional Layer Stages

A typical setup of a convolutional layer is as follows:

1. Convolution stage: affine transform

2. Detector stage: nonlinearity (e.g., ReLU)

3. Pooling stage: locally combine activities in some way (max, avg, ...)

Locality of the item that activated the neurons isn't too important, so we profit from dimensionality reduction. Another thing that turns out to be so is that most of the kernels that are learned resemble a low-pass filter. Hence, when we sub-sample the images most of the information is still contained.

12.9 — Pooling

The most frequently used pooling function is: *max pooling*. But one can imagine using other pooling functions such as min, avg, and softmax.

D. (Max-Pooling)

Max pooling works, as follows, if we define a window size of $r=3$ (in 1D or 2D), then

- $1D: x_{ij}^{\text{max}} = \max\{x_{i+k} \mid 0 \leq k < r\}$
- $2D: x_{ij}^{\text{max}} = \max\{x_{i+k, j+l} \mid 0 \leq k, l < r\}$

So, in general we just take the maximum over a small "patch"/"neighbourhood" of some units.

T. (Max-Pooling: Invariance)

Let \mathcal{T} be the set of invertible transformations (e.g., integral transforms, integral operators). Then \mathcal{T} forms a group w.r.t. function composition: $(\mathcal{T}, \circ, \text{id})$.

12.10 — Sub-Sampling (aka "Strides")

Often, it is desirable to reduce the size of the feature maps. That's why sub-sampling was introduced.

D. (Sub-Sampling)

Hereby the temporal/spatial resolution is reduced.

Com. Often, the sub-sampling is done via a max-pooling according to some interval step size (a.k.a. stride)

- Loss of information

+ Dimensionality reduction

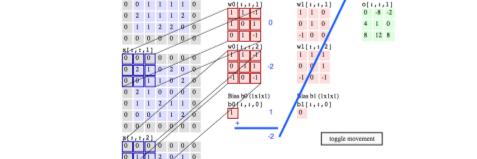
+ Increase of efficiency

12.11 — Channels

Ex. Here we have

- an input signal that is 2D with 3 channels ($7 \times 7 \times 3$) (image x channels)

- and we want to learn two filters W_0 and W_1 , which each process the 3 channels, and sum the results of the convolutions across each channel leading to a tensor of size $3 \times 3 \times 2$ (convolution result x num convolutions)



Usually we convolve over all of the channels together, such that each convolution has the information of all channels at its disposition and the order of the channels hence doesn't matter.

12.12 — CNNs in Computer Vision

So the typical use of convolution that we have in vision is: a sequence of convolutions

1. that *reduce* the spatial dimensions (sub-sampling)

2. that *increase* the number of channels

The deeper we go in the network, we transform the spatial information into a semantic representation. Usually, most of the parameters lie in the fully connected layers

12.13 — Famous CNN Architectures

12.13.1 — LeNet, 1989

MNIST, 2 Convolutional Layers + 2 Fully-connected layers

12.13.2 — AlexNet

MNIST, 3 Convolutional Layers (with max-pool subsampling) + 1 Fully connected layer

12.13.3 — AlexNet

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

As we can see, the convergence time is bounded by a time that depends on our initial guess, and the Lipschitz constant L .

However, things become even stranger because of the curvature. As we can see, the gradient norm gets larger and larger the more we train (can be checked empirically). And the gradient norm also tends to have larger fluctuations. And as we can see, starting at some point, the error just fluctuates around at a certain level. Actually, one might assume that as we're getting closer to the minimum, the gradient should get smaller and smaller, as the objective gets flatter and flatter at the optimal point - but that's actually not the case!

Com. So we have a polynomial (linear) convergence rate of θ towards the optimal parameter θ^* (note: just in the convex setting!). As we can see, the convergence time is bounded by a time that depends on our initial guess, and the Lipschitz constant L .

Now, taking the derivatives w.r.t. the parameters, we get (using the chain rule)

$$\nabla_Q \mathcal{R}(Q, W) = \frac{\partial \mathcal{R}(Q)}{\partial Q} \frac{\partial A}{\partial Q}$$

Now, from what we've seen before, we can express the risk as (due to trace identities, trace linearity, etc.) just by replacing $A = QW$, so

$$\mathcal{R}(Q, W) = \text{const.} + \text{Tr}((QW)(QW)^T) - 2\text{Tr}(QW^T)$$

Now, taking the derivatives w.r.t. the parameters, we get (using the chain rule)

$$\nabla_Q \mathcal{R}(Q, W) = \frac{\partial \mathcal{R}(Q)}{\partial Q} \frac{\partial A}{\partial Q}$$

So since μ and σ are functions of the weights and they are differentiated.

A further note about batch-normalization is that it doesn't change the information in the data, because since we have μ and σ we could theoretically recuperate the original activations.

Now, some implementation details:

- The bias term before the batch normalization should be removed (since we're removing the mean it makes no sense).

- At training time, the statistics are computed per batch, hence they're very noisy. So what people do in practice (e.g., when they're predicting just one sample) is that they keep a running average over the batch batch-norm statistics. So, at test time, μ and σ are replaced by the running averages that were collected during training time. An alternative, is to pass through the whole dataset at the end of the training and re-compute the statistics - that may work even better (but it takes a lot of time).

What's not very clear is why batch-normalization works. The original paper about batch-normalization (BN) said that BN reduces the variance of the features. The reason for this is that BN shifts the mean to zero. So BN is a very simple classifier, that will basically classify everything that is negative to one class, and everything that is positive to another class. Then, when we shift the data by a constant vector, then, without batch-normalization we'd shift all the datapoints into one class. However, with BN the mean is removed we'll remove that constant shift the BN layer and it will work out. So BN reduces the covariance shift. That was the effect that the inventors of BN described.

However, it turns out that some other people came later on and said the following: They say that the effect of the covariance-shift is that it makes the landscape of the loss more smooth. Hence, the optimization works better and gives better results.

However, no one really knows why BN works so well.

13.8.8 — Other Heuristics

Curriculum learning and non-uniform sampling

of data points → focus on the most relevant examples (Bengio, Laradour, Collobert, Weston, 2009) (DL-Book: 8.7.6) Or

increase hardness of tasks (corner-cases) as BN improves

Continuation methods

: define a family of (smooth) objective

functions and track solutions, gradually change hardness of loss (DL-Book: 8.7.6)

Heuristics for initialization

(DL-Book: 8.4) scale the weights

of each layer in a way that at the end of the layer, the data has more or less the same energy (and gradient norms are more or less the same at each layer).

Pre-training

(DL-Book: 8.7.4). for better initialization, to avoid local minima (less relevant today).

13.10 — Norm-Based Regularization

$\mathcal{R}_\Omega(Q; S) = \mathcal{R}(Q; S) + \Omega(Q)$,

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)

$\Omega(Q) = \frac{1}{2} \sum_{i=1}^L \lambda_i^2 \|W_i\|_F^2, \quad \lambda_i^2 \geq 0$

Com. It's common practice to only penalize the weights, and not the biases.

The assumption here is that the weights have to be small.

Then at every timestep $t = 1, 2, \dots$ we pick a random subset

$$S_t \subseteq S_N, \quad K \leq N. \quad (\text{usually } K \ll N)$$

And since we're picking S_t at random, note how

$$\mathbb{E}[\mathcal{R}(S_t)] = \mathcal{R}(S_N)$$

And thus it also holds for the gradient that

$$\nabla_\theta \mathcal{R}(S_t) \stackrel{\text{lin.}}{\approx} \nabla_\theta \mathcal{R}(S_N) = \nabla_\theta \mathcal{R}(S_N).$$

So, with SGD, we just do gradient descent, where at each t we'll do a randomization of S_t , so

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

Initialization: $\theta(0) = \mathbf{o}$, $v(0) = \mathbf{o}$

$$\mathbb{E}[V] = \mathbb{E}[-\nabla_\theta \mathcal{R}]$$

Typical values: $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$

Then at every timestep $t = 1, 2, \dots$ we pick a random subset

$$g(t) = \nabla_\theta \mathcal{R}(\theta(t-1)) \quad (\text{gradient})$$

$m(t) = \beta_1 m(t-1) + (1-\beta_1) g(t)$ (update biased first moment estimate)

$v(t) = \beta_2 v(t-1) + (1-\beta_2) g(t)^2$ (update b. second raw moment estimate)

$\bar{m}(t) = m(t)/(1-\beta_1^t)$ (bias correction first moment estimate)

$\bar{v}(t) = v(t)/(1-\beta_2^t)$ (bias correction second raw moment estimate)

$\theta(t) = \theta(t-1) + \eta \bar{m}(t)/(\sqrt{\bar{v}(t)} + \epsilon)$ (update params)

13.9.7 — Batch Normalization

Batch normalization (Ioffe & Szegedy, 2015) is one of the most controversial but most useful tricks in DL.

One of the big problems that we have when we optimize NNs is that usually there exist

from right to left (starting at w^*) the risk has a very low eigenvalue, and hence \hat{w} is reduced much more along that dimension.

D. (L1-Regularization (sparsity inducing))

$$\Omega(\theta) = \sum_{l=1}^L \lambda^l \| \mathbf{W}^l \|_1 = \sum_{l=1}^L \lambda^l \sum_{i,j} |w_{ij}|, \quad \lambda^l \geq 0$$

- 13.10.1—Regularization via Constrained Optimization

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

$$\min_{\theta, \|\theta\| \leq r} \mathcal{R}(\theta).$$

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

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

where

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

So we're essentially clipping the weights.

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

Hinton made some research in 2013 and realized that

- the constraints do not affect the initial learning (as the weights are assumed to be small at the beginning), so we won't clip the weights.
- So the constraints only become active, once the weights are large.

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

This had some practical success in stabilizing the optimization.

- 13.10.2—Early Stopping

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

- 14.1—Natural Language Processing

Similarities between text and image processing: local information. Differences between text and images: texts have various lengths, texts may have long-term interactions, language is a man-made conception on how to communicate with each other / pictures capture the reality, pictures capture the reality / sentences may mean different things in different contexts

- 14.1.1—Word Embeddings

Basic Idea: Map symbols over a vocabulary V to a vector representation = embedding into an (euclidean) vector space (see lookup table in architecture overview).

embedding map: (vocabulary) $V \mapsto \mathbb{R}^d$ (embeddings)

(symbolic) $w \mapsto \mathbf{x}_w$ (quantitative)

word $w \in V \rightarrow$ one-hot $w \in \{0, 1\}^{|V|} \rightarrow$ embedding \mathbf{x}_w .

$m := |V|$, usually $|V| = 10^5$

d = dimensionality of embedding, $d \ll m$

So for each of the m words in V we have a corresponding embedding in \mathbb{R}^d , which can be stored in a shared lookup table:

$\mathcal{R}^{d \times m}$ shared lookup table

Any sentence of k words can then be represented as a $d \times k$ matrix (a sequence of k embedding vectors in \mathbb{R}^d).

Now, how should an embedding be? Ideally, the embedding carries the information/structure that we need in order to go from the input text to the question that we want to solve. Typical questions are:

- Clustering based on context (co-occurrence)
- Sentiment analysis (group words according to mood/feelings)
- Translation (group by meaning)
- Part-of-Speech tagging (understand the structure of text, e.g., location, time, actor, ... or, noun, verb, adjective, ...)

- 14.1.1.2—Bi-Linear Models

The first thing that we could do is to use an information theoretic quantity: the so-called *mutual information*. The mutual information is described in information theory as *how much information on random variable has about another random variable*. If two variables are independent, then, the mutual information will be zero.

So, if we put two words nearby, it's because they have to be related somehow in the *meaning* of the sentence. Hence, we expect them to have a larger mutual information.

D. (Pointwise Mutual Information)

pmi(v, w) = $\log \left(\frac{P(v, w)}{P(v) P(w)} \right) = \log \left(\frac{P(v|w)}{P(v)} \right) \approx \mathbf{x}_v^\top \mathbf{x}_w + \text{const}$

Com. As you can see this metric is bi-linear.

We interpret the vectors as *latent variables* and link them to the observable probabilistic model. So the pointwise mutual information is related to the inner product between the latent vectors (the more related, the more co-linear the latent representations have to be).

- 13.11.1—Invariant Architectures

Instead of augmenting the dataset one could build an architecture that is invariant to certain transformations of the data.

First, we distinguish the following terms: Let's say we have some x and apply the transformation $\mathbf{x} := \tau(\mathbf{x})$. Then for our neural network F :

- **D. (Invariance)** means that $F(\mathbf{x}) = F(\tau(\mathbf{x}))$.
- **Equivariance** means that $\tau(F(\mathbf{x})) = F(\tau(\mathbf{x}))$.

So applying the transformation before or after applying F doesn't change a thing (e.g., convolutions and translations are equivariant).

D. (Co-Occurrence Set)

Here we look at the *pairwise occurrences* of words in a *context window* of size R . So, if we have a long sequence of words $\mathbf{w} = (w_1, \dots, w_T)$, then the co-occurrence index set is defined as

$$C := \{(i, j) \mid 1 \leq i < j \leq R\}.$$

D. (Co-Occurrence Matrix)

Note that in order to get an (empirical) idea of the co-occurrence frequencies one could compute the co-occurrence matrix

$$\log(P(w_t \mid \mathbf{w} := w_{t-1}, \dots, w_1)) = \mathbf{x}_{w_t}^\top \mathbf{z}_w + \text{const}$$

Properties: $\mathbf{C} = \mathbf{C}^\top$ (symmetric), peaky, sparse.

One approach for embeddings: do PCA of \mathbf{C} and use k eigenvectors corresponding to largest eigenvalues of \mathbf{C} . Note that we have

$$\begin{aligned} C_{ij} &= \text{one-hot}(w_i) \mathbf{C} \text{one-hot}(w_j) = \mathbf{o}_i \mathbf{V} \mathbf{A}^\top \mathbf{o}_j \\ &\approx \mathbf{o}_i \mathbf{V} \mathbf{k}_k \mathbf{V}^\top \mathbf{o}_j = \mathbf{o}_i \mathbf{V} \mathbf{k}_k \frac{1}{k} \mathbf{k}_k^\top \mathbf{V} \mathbf{o}_j \end{aligned}$$

- 14.2.1—ConvNets: Word Representations

de-embedding: $\mathbf{V} \mathbf{A}^{-\frac{1}{2}}$ (then find nearest neighbour)

Problem: \mathbf{C} is huge ($|V|^2$), hence matrix-factorization becomes prohibitively expensive!

Solution: Use skip-gram approach to avoid computing \mathbf{C} at all!

The solution to this is pretty simple: we train a model that tries to predict for one word w the preceding and following words:

$$w_{t-1}, w_{t-1+1}, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_{t+c-1}, w_{t+c}.$$

Here's an illustration of the model for $t = 3$: input $w(t)$ → projection → $w(t-2), w(t-1), w(t+1), w(t+2)$ output

Note that the assumption (or simplification) of this model is that it assumes that the words w_i, w_j within the window $+c, -c$ of w_t are conditionally independent of each other given w_t .

$$W_t \perp W_j \mid W_i \quad (i \neq j \wedge i \neq t \wedge j \neq t)$$

That might be too much of an assumption but you can see that sometimes when we're talking about something, we may change the order of the words and still mean the same thing (e.g., "I was born in 1973" → "1973 is the year I was born"). So in a way, we're just trying to capture the meaning of w_t with this. So this gives us an idea of the context of w_t and might relieve the structure we're looking for. So, it's not as optimal as computing \mathbf{C} , but it's a way to start.

Benefits: benefits of ensembles with the runtime complexity of the training of one network. The network gets trained to have many different paths through it to get the right result (as neurons are turned off).

Equivalent to: adding multiplicative noise to weights or training exponentially many sub-networks $\sum_{n=1}^N \binom{N}{n} = 2^N$ where n is the number of compute units (so at each iteration we turn some nodes off according to some probability). So we're getting the benefits of ensemble with the runtime complexity of just training one network.

Ensembling corresponds to taking geometric mean (instead of usual arithmetic) (must have to do with exponential growth of networks) of the ensembles:

$$\text{Ensemble } (\mathbf{y} \mid \mathbf{x}) = \sqrt[d]{\prod_{\mu} P(\mu) P(y \mid \mathbf{x}, \mu)}$$

Having to sample several sub-networks for a prediction is somewhat inconvenient, so the idea that Hinton et al. came up with is: scaling each weight w_{ij}^t by the probability of the unit j being active:

$$\begin{aligned} w_{ij}^t &\leftarrow p_{j-1}^t w_{ij}^t \\ &= \arg \max_{\theta} \sum_{(i,j) \in C_R} \log \left(\frac{\exp(\mathbf{x}_w^\top \mathbf{z}_{w_j})}{\sum_{u \in V} \exp(\mathbf{x}_w^\top \mathbf{z}_{w_u})} \right) \end{aligned}$$

This makes sure that the net (total) input to unit x_j^t is calibrated, i.e.,

$$\sum_j w_{ij}^t p_{j-1}^t = \mathbb{E}_{Z \sim P(Z)} \left[\sum_j p_{j-1}^t w_{ij}^t x_j^t \right] = \sum_j p_{j-1}^t w_{ij}^t x_j^t$$

It can be shown that this approach leads to a (sometimes exact) approximation of a geometrically averaged ensemble (see DL-Book, 7.12).

Ex. Let's say that at the end we selected each unit with a probability of 0.5. Then when typically when we're finished with training our neural network, we're going to multiply all the weights that we obtained with 0.5 to reduce the contribution of each of the features (since we'll have all of them). So with this trick for the prediction we can just do a single forward pass.

- 14.3—Natural Language Processing

Similarities between text and image processing: local information. Differences between text and images: texts have various lengths, texts may have long-term interactions, language is a man-made conception on how to communicate with each other / pictures capture the reality, pictures capture the reality / sentences may mean different things in different contexts

- 14.3.1—Word Embeddings

Basic Idea: Map symbols over a vocabulary V to a vector representation = embedding into an (euclidean) vector space (see lookup table in architecture overview).

embedding map: (vocabulary) $V \mapsto \mathbb{R}^d$ (embeddings)

(symbolic) $w \mapsto \mathbf{x}_w$ (quantitative)

word $w \in V \rightarrow$ one-hot $w \in \{0, 1\}^{|V|} \rightarrow$ embedding \mathbf{x}_w .

$m := |V|$, usually $|V| = 10^5$

d = dimensionality of embedding, $d \ll m$

So for each of the m words in V we have a corresponding embedding in \mathbb{R}^d , which can be stored in a shared lookup table:

$\mathcal{R}^{d \times m}$ shared lookup table

Any sentence of k words can then be represented as a $d \times k$ matrix (a sequence of k embedding vectors in \mathbb{R}^d).

Now, how should an embedding be? Ideally, the embedding carries the information/structure that we need in order to go from the input text to the question that we want to solve. Typical questions are:

- Clustering based on context (co-occurrence)
- Sentiment analysis (group words according to mood/feelings)
- Translation (group by meaning)
- Part-of-Speech tagging (understand the structure of text, e.g., location, time, actor, ... or, noun, verb, adjective, ...)

- 14.3.1.2—Bi-Linear Models

The first thing that we could do is to use an information theoretic quantity: the so-called *mutual information*. The mutual information is described in information theory as *how much information on random variable has about another random variable*. If two variables are independent, then, the mutual information will be zero.

So, if we put two words nearby, it's because they have to be related somehow in the *meaning* of the sentence. Hence, we expect them to have a larger mutual information.

D. (Pointwise Mutual Information)

$$\text{pmi}(v, w) = \log \left(\frac{P(v, w)}{P(v) P(w)} \right) = \log \left(\frac{P(v|w)}{P(v)} \right) \approx \mathbf{x}_v^\top \mathbf{x}_w + \text{const}$$

Com. As you can see this metric is bi-linear.

We interpret the vectors as *latent variables* and link them to the observable probabilistic model. So the pointwise mutual information is related to the inner product between the latent vectors (the more related, the more co-linear the latent representations have to be).

- 13.11.1—Invariant Architectures

Instead of augmenting the dataset one could build an architecture that is invariant to certain transformations of the data.

First, we distinguish the following terms: Let's say we have some x and apply the transformation $\mathbf{x} := \tau(\mathbf{x})$. Then for our neural network F :

- **D. (Invariance)** means that $F(\mathbf{x}) = F(\tau(\mathbf{x}))$.
- **Equivariance** means that $\tau(F(\mathbf{x})) = F(\tau(\mathbf{x}))$.

So applying the transformation before or after applying F doesn't change a thing (e.g., convolutions and translations are equivariant).

D. (Co-Occurrence Set)

Here we look at the *pairwise occurrences* of words in a *context window* of size R . So, if we have a long sequence of words $\mathbf{w} = (w_1, \dots, w_T)$, then the co-occurrence index set is defined as

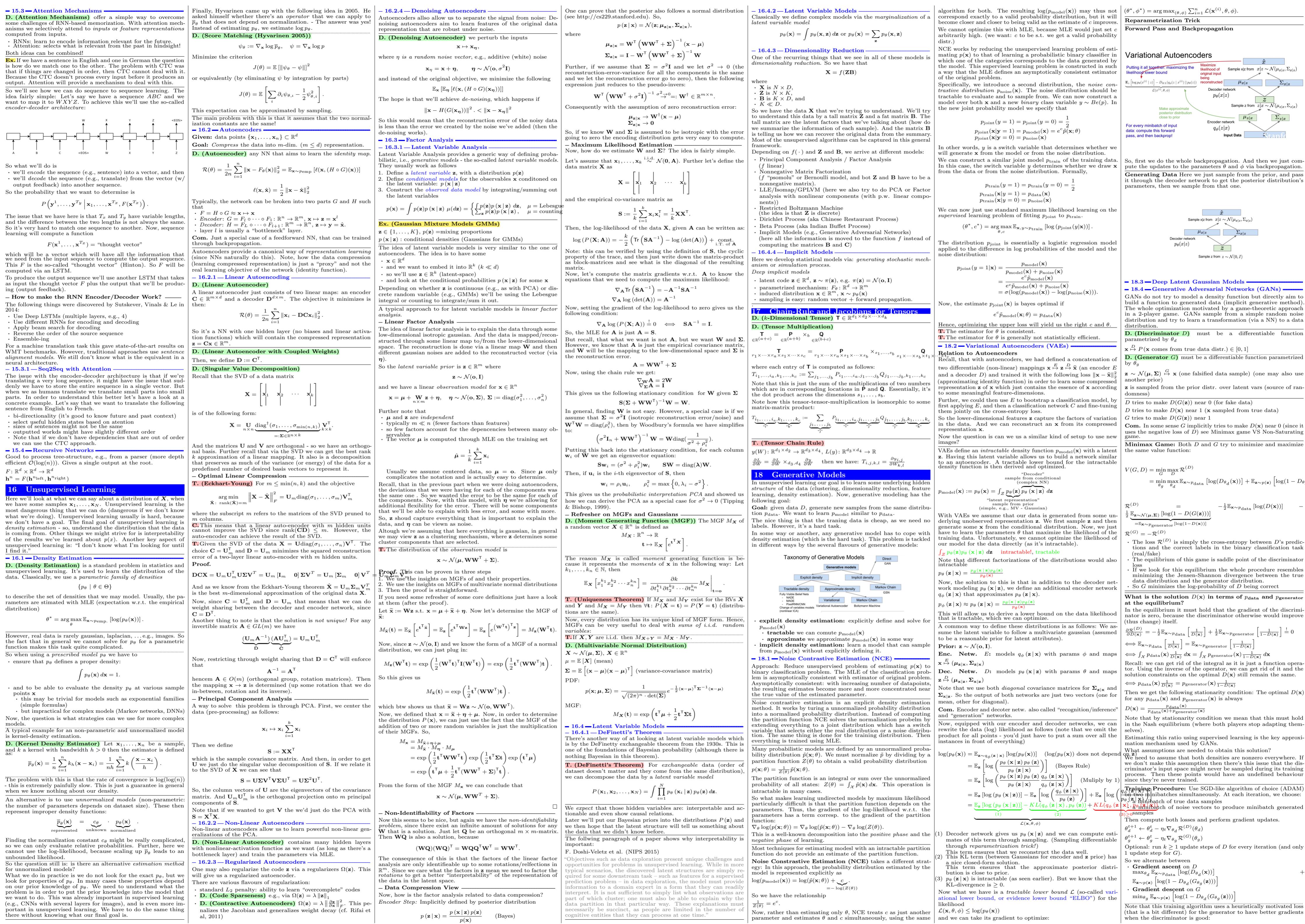
$$C := \{(i, j) \mid 1 \leq i < j \leq R\}.$$

D. (Co-Occurrence Matrix)

Note that in order to get an (empirical) idea of the co-occurrence frequencies one could compute the co-occurrence matrix

$$\log(P(w_t \mid \mathbf{w} := w_{t-1}, \dots, w_1)) = \mathbf{x}_{w_t}^\top \mathbf{z}_w + \text{const}$$

Properties: $\mathbf{C} = \mathbf{C}^\top$ (symmetric), peaky, sparse.



```

for number of training iterations do
    for k steps do
        • Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
        • Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{\text{data}}(x)$ .
        • Update the discriminator by ascending its stochastic gradient:
            
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

    end for
    • Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
    • Update the generator by ascending its stochastic gradient (improved objective):
        
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

end for

```