

**1 Probability**

Sum Rule  $P(X = x_i) = \sum_{j=1}^n P(X = x_i, Y = y_j)$   
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.  $X \perp\!\!\!\perp Y|Z \Rightarrow P(X|Y|Z) = P(X|Z)$

$E[X] = \int_X t \cdot f_X(t) dt = \mu_X$   
 $\text{Var}[X] = E[(X - E[X])^2] = \int_X (t - E[X])^2 f_X(t) dt = E[X^2] - E[X]^2$   
 $\text{Cov}(X, Y) = E_{x,y}[(X - E[X])(Y - E[Y])]$   
 $\text{Cov}(X) = \text{Cov}(X, X) = \text{Var}[X]$   
 $X, Y \text{ independent} \Rightarrow \text{Cov}(X, Y) = 0$

$\text{Var}[\mathbf{A}\mathbf{X}] = \mathbf{A}\text{Var}[\mathbf{X}]\mathbf{A}^\top$   $\text{Var}[a\mathbf{X} + b] = a^2 \text{Var}[\mathbf{X}]$   
 $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + 2 \sum_{i,j,i < j} a_i a_j \text{Cov}(X_i, X_j)$   
 $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + \sum_{i,j,i \neq j} a_i a_j \text{Cov}(X_i, X_j)$   
 $\frac{\partial}{\partial t} \mathbf{F}(t) = \mathbf{F}'(t)$  ( $\mathbf{F}$  derivative of c.d.f.)  
 $f_{\mathbf{A}}(y) = \frac{1}{a} f_Y(\frac{y}{a})$   
Empirical CDF:  $\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq t\}}$   
Empirical PDF:  $\hat{f}_n(t) = \frac{1}{n} \sum_{i=1}^n \delta(t - X_i)$  (continuous)  
Empirical PDF:  $\hat{p}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{t \in D\}}$  (discrete)

**2.4.4 Scalar-by-Matrix**

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{Xb}] = \mathbf{ab}^\top$   $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{Xb}] = \mathbf{X}(\mathbf{ab}^\top + \mathbf{ba}^\top)$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{ba}^\top$   $\frac{\partial}{\partial \mathbf{x}} [\text{Tr}(\mathbf{X})] = \mathbf{I}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{Xa}] = \frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{a}] = \frac{\partial}{\partial \mathbf{x}} [\text{Tr}(\mathbf{AXB})] = \mathbf{A}^\top \mathbf{B}^\top$   
 $\mathbf{aa}^\top$   $\frac{\partial}{\partial \mathbf{x}} [\text{Tr}(\mathbf{AX}^\top \mathbf{B})] = \mathbf{BA}$

**4.5 Vector-by-Matrix (Generalized Gradient)**

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{Xa}] = \mathbf{X}^\top$

**5 General Machine Learning**

**Likelihood**  $H := \{x \mid \langle w, x - p \rangle = 0\} = \{x \mid \langle w, x \rangle = b\}$   
where  $w = (w, p)$ ,  $w$  = normal vector,  $p$  points onto a point on the plane.

**D. (Level Sets)** of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a one-parametric family of sets defined as  
 $L_f(c) := \{x \mid f(x) = c\} = f^{-1}(c) \subseteq \mathbb{R}^n$ .

**6 Information Theory**

**D. (Entropy)** Let  $X$  be a random variable distributed according to  $p(X)$ . Then the entropy of  $X$  is  
 $H(X) = -\sum_{x \in \mathcal{X}} p(x) \log(p(x)) = E[I(X)] = E[-\log(P(X))]$  ≥ 0. describes the expected information content  $I(X)$  of  $X$ .

**D. (Cross-Entropy)** for the distributions  $p$  and  $q$  over a given set is  
 $H(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(q(x)) = E_{x \sim p}[-\log(q(x))]$  ≥ 0.  
 $H(X; p, q) = H(X) + K L(p, q) \geq 0$ , where  $H$  uses  $p$ .  
 $\text{Com.}$  The second formulation clearly shows why  $q := p$  is the minimizer of the cross-entropy (or hence: the maximizer of the likelihood).

**Com.** Usually,  $q$  is the approximation of the unknown  $p$ .  
**Relation to Log-Likelihood** In classification problems we want to estimate the probability of different outcomes. If we have the following quantities:  
- estimated probability of outcome  $i$  is  $q_i$ . Now we want to tune  $q$  in a way that the data gets the most likely. First, let's just see how good  $q$  is doing.  
- the frequency (empirical probability) of outcome  $i$  in the data is  $p_i$ .  
-  $n$  data points  
Then the likelihood of the data under  $p_i$  is  
 $\prod_{i=1}^n q_i^{n-p_i}$  since the model estimates event  $i$  with probability  $q_i$  exactly  $n$  not  $p_i$  times. Now the log-likelihood, divided by  $n$  is  
 $\frac{1}{n} \sum_{i=1}^n n p_i \log(q_i) = \sum_{i=1}^n p_i \log(q_i) = -H(p, q)$  Hence, maximizing the log-likelihood corresponds to minimizing the cross-entropy (which is why it's used so often as a loss).

**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  
 $KL(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log\left(\frac{q(x)}{p(x)}\right) = \sum_{x \in \mathcal{X}} p(x) \log\left(\frac{p(x)}{q(x)}\right) \geq 0$   
 $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$   $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x - \hat{\mu})(x - \hat{\mu})^\top$

**T. (Chebyshev)** Let  $X$  be a rv with  $E[X] = \mu$  and variance  $\text{Var}[X] = \sigma^2 < \infty$ . Then for any  $\epsilon > 0$ , we have  $P(|X - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$ .

**2 Analysis**

Log-Trick (Identity):  $\nabla_\theta [p_\theta(\mathbf{x})] = p_\theta(\mathbf{x}) \nabla_\theta [\log(p_\theta(\mathbf{x}))]$

**T. (Cauchy-Schwarz)**  
 $\mathbf{u}, \mathbf{v} \in V: (\mathbf{u}, \mathbf{v}) \leq \|(\mathbf{u}, \mathbf{v})\| \leq \|\mathbf{u}\| \|\mathbf{v}\|$ .  
 $\mathbf{u}, \mathbf{v} \in V: 0 \leq (\mathbf{u}, \mathbf{v}) \leq \|\mathbf{u}\| \|\mathbf{v}\|$ .  
Special case:  $(\sum_i x_i y_i)^2 \leq (\sum_i x_i^2)(\sum_i y_i^2)$ .  
Special case:  $E[XY]^2 \leq [E[X]^2][E[Y]^2]$ .

**T. (Fundamental Theorem of Calculus)**  
 $f(y) - f(x) = \int_{[x,y]} \nabla f(\tau) \cdot d\tau = \int_{t=0}^1 \nabla f(\gamma(t))^\top \gamma'(t) dt$   
 $f(y) - f(x) = \int_0^1 \nabla f((1-t)x + ty) \cdot (y - x) dt$   
 $\text{Com.}$  Create a path  $\gamma$  from  $x$  to  $y$  and integrate the dot product of the gradient of the function-values at the path with the derivative of the path.

**D. (Saddle Points etc.)**  
**T. (Jensen)**  $f$  convex/concave:  $\forall i: \lambda_i \geq 0, \sum_{i=1}^n \lambda_i = 1$   
 $f(\sum_{i=1}^n \lambda_i x_i) \leq \sum_{i=1}^n \lambda_i f(x_i)$   
Special case:  $f(E[X]) \leq E[f(X)]$ .

**D. (Lagrangian Formulation)** of  $f(x, y)$  s.t.  $g(x, y) = c$   
 $L(x, y, \gamma) = f(x, y) - \gamma(g(x, y) - c)$

**3 Linear Algebra**

**T. (Positive Criterion)** A  $d \times d$  matrix is positive semi-definite if and only if all the upper left  $k \times k$  for  $k = 1, \dots, d$  have a positive determinant.  
Negative definite:  $\det < 0$  for all odd-sized minors, and  $\det > 0$  for all even-sized minors.  
Otherwise: indefinite.

**D. (Trace)** of  $A \in \mathbb{R}^{n \times n}$  is  $\text{Tr}(A) = \sum_{i=1}^n a_{ii}$ .

**4 Derivatives**

**4.1 Numerator and Denominator Convention**  
**Jacobian Layout Convention** We use the denominator-layout. For a vector-valued function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  we define  
 $\left(\frac{\partial f}{\partial \mathbf{x}}\right)_{ij} := \frac{\partial f_j}{\partial x_i}, \quad \frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{n \times m}$ .  
Hence, gradients of scalar-valued functions are column vectors, and the chain rule takes the form  
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{g}(\mathbf{x}))] = \frac{\partial \mathbf{f}}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{x}}$ .

**Remark.** Sometimes the numerator-layout is used, where the Jacobian is defined as  $(\mathbf{J}_x)_{ij} = \frac{\partial f_i}{\partial x_j} \in \mathbb{R}^{m \times n}$ . The two conventions are related by transposition.

**4.2 Scalar-by-Vector**

**Denominator Convention**  $\frac{\partial}{\partial \mathbf{x}} [\mathbf{u}(\mathbf{x})\mathbf{v}(\mathbf{x})] = \mathbf{u}(\mathbf{x}) \frac{\partial \mathbf{v}(\mathbf{x})}{\partial \mathbf{x}} + \mathbf{v}(\mathbf{x}) \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{u}(\mathbf{x})\mathbf{g}(\mathbf{x})] = \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}} \mathbf{g}(\mathbf{x}) + \mathbf{u}(\mathbf{x}) \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})\mathbf{Ag}(\mathbf{x})] = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \mathbf{Ag}(\mathbf{x}) + \mathbf{A}^\top \mathbf{f}(\mathbf{x})\mathbf{a}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{x}^\top \mathbf{x}] = 2\mathbf{x}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{b}^\top \mathbf{Ax}] = \mathbf{A}^\top \mathbf{b}$   
 $\frac{\partial}{\partial \mathbf{x}} [(\mathbf{Ax} + \mathbf{b})^\top \mathbf{C}(\mathbf{Dx} + \mathbf{e})] = \mathbf{D}^\top \mathbf{C}^\top (\mathbf{Ax} + \mathbf{b}) + \mathbf{A}^\top \mathbf{C}(\mathbf{Dx} + \mathbf{e})$   
 $\frac{\partial}{\partial \mathbf{x}} [\| \mathbf{f}(\mathbf{x}) \|^2] = \frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})^\top \mathbf{f}(\mathbf{x})] = 2 \frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2 \mathbf{J}_f(\mathbf{x})$

**4.3 Vector-by-Vector**

**A, C, D, a, b, e** not a function of  $\mathbf{x}$ ,  
 $\mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{x}), \mathbf{h}(\mathbf{x}), \mathbf{u} = \mathbf{u}(\mathbf{x}), \mathbf{v} = \mathbf{v}(\mathbf{x})$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{u}(\mathbf{x})\mathbf{f}(\mathbf{x})] = \mathbf{u}(\mathbf{x}) \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} + \mathbf{f}(\mathbf{x}) \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{x} \otimes \mathbf{a}] = \text{diag}(\mathbf{a})$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a} \otimes \mathbf{b}] = \mathbf{0}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{x}] = \mathbf{I}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{Ax}] = \mathbf{A}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{x}^\top \mathbf{A}] = \mathbf{A}^\top$

**4.4.4 Scalar-by-Matrix**

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{Xb}] = \mathbf{ab}^\top$   $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{Xb}] = \mathbf{X}(\mathbf{ab}^\top + \mathbf{ba}^\top)$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{ba}^\top$   $\frac{\partial}{\partial \mathbf{x}} [\text{Tr}(\mathbf{X})] = \mathbf{I}$   
 $\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{Xa}] = \frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^\top \mathbf{X}^\top \mathbf{a}] = \frac{\partial}{\partial \mathbf{x}} [\text{Tr}(\mathbf{AXB})] = \mathbf{A}^\top \mathbf{B}^\top$   
 $\mathbf{aa}^\top$   $\frac{\partial}{\partial \mathbf{x}} [\text{Tr}(\mathbf{AX}^\top \mathbf{B})] = \mathbf{BA}$

**4.5 Vector-by-Matrix (Generalized Gradient)**

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{Xa}] = \mathbf{X}^\top$

**5 General Machine Learning**

**Likelihood**  $H := \{x \mid \langle w, x - p \rangle = 0\} = \{x \mid \langle w, x \rangle = b\}$   
where  $w = (w, p)$ ,  $w$  = normal vector,  $p$  points onto a point on the plane.

**D. (Level Sets)** of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a one-parametric family of sets defined as  
 $L_f(c) := \{x \mid f(x) = c\} = f^{-1}(c) \subseteq \mathbb{R}^n$ .

**6.1 Information Theory**

**D. (Entropy)** Let  $X$  be a random variable distributed according to  $p(X)$ . Then the entropy of  $X$  is  
 $H(X) = -\sum_{x \in \mathcal{X}} p(x) \log(p(x)) = E[I(X)] = E[-\log(P(X))]$  ≥ 0. describes the expected information content  $I(X)$  of  $X$ .

**T. (Comp. of Lin. Maps/- is a Lin. Map/Unit)**  
Let  $F_1, \dots, F_L$  be linear maps, then  $F = F_L \circ \dots \circ F_2 \circ F_1$  is also a linear map.  
**C.** Every L-layer NN of linear layer collapses to a 1-layer NN. Further note that hereby  $\text{rank}(F) \leq \min_{i \in \{1, \dots, L\}} \text{rank}(F_i)$ .

**D. (Cross-Entropy)** for the distributions  $p$  and  $q$  over a given set is  
 $H(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(q(x)) = E_{x \sim p}[-\log(q(x))]$  ≥ 0.  
 $H(X; p, q) = H(X) + K L(p, q) \geq 0$ , where  $H$  uses  $p$ .  
 $\text{Com.}$  The second formulation clearly shows why  $q := p$  is the minimizer of the cross-entropy (or hence: the maximizer of the likelihood).

**Com.** Usually,  $q$  is the approximation of the unknown  $p$ .  
**Relation to Log-Likelihood** In classification problems we want to estimate the probability of different outcomes. If we have the following quantities:  
- estimated probability of outcome  $i$  is  $q_i$ . Now we want to tune  $q$  in a way that the data gets the most likely. First, let's just see how good  $q$  is doing.  
- the frequency (empirical probability) of outcome  $i$  in the data is  $p_i$ .  
-  $n$  data points  
Then the likelihood of the data under  $p_i$  is  
 $\prod_{i=1}^n q_i^{n-p_i}$  since the model estimates event  $i$  with probability  $q_i$  exactly  $n$  not  $p_i$  times. Now the log-likelihood, divided by  $n$  is  
 $\frac{1}{n} \sum_{i=1}^n n p_i \log(q_i) = \sum_{i=1}^n p_i \log(q_i) = -H(p, q)$  Hence, maximizing the log-likelihood corresponds to minimizing the cross-entropy (which is why it's used so often as a loss).

**D. (Ridge Function)**  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a ridge function, if it can be written as  $f(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b)$  for some  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ ,  $\mathbf{w} \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ .  
**Rep.** polygonal lines by (linear function +) linear combinations of  $(-)$ - or  $(+)$ -functions

**3. Apply dimension lifting lemma to show density of the linear span of resulting ridge function families  $\mathcal{G}_{++}^n$  and  $\mathcal{G}_{-+}^n$ .**

**11.5.1 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):  
 $H^+ := \{x \mid \mathbf{w}^\top \mathbf{x} + b > 0\} \subseteq \mathbb{R}^n$   
 $H^- := \{x \mid \mathbf{w}^\top \mathbf{x} + b < 0\} \subseteq \mathbb{R}^n$   
 $H^0 := \{x \mid \mathbf{w}^\top \mathbf{x} + b = 0\} = \mathbb{R}^n - H^+ - H^- \subseteq \mathbb{R}^n$

**D. (Loss Function)**  
In decision theory, we strive to minimize the expected risk (defined through a loss function  $\ell$ ) of a function  $F$ .

**D. (Expected Risk of a Function F)**  
 $F^* = \arg \min_F \sum_{y \in \mathcal{Y}} \int_{\mathcal{X}} \ell(y, F(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x}$   
 $= \arg \min_F \mathbb{E}_{\mathbf{x} \sim \mathcal{X}} [\ell(Y, F(\mathbf{x}))]$   
 $\mathcal{R}^*(F)$  expected risk of  $F$

**11.5.2 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):  
 $H^+ := \{x \mid \mathbf{w}^\top \mathbf{x} + b > 0\} \subseteq \mathbb{R}^n$   
 $H^- := \{x \mid \mathbf{w}^\top \mathbf{x} + b < 0\} \subseteq \mathbb{R}^n$   
 $H^0 := \{x \mid \mathbf{w}^\top \mathbf{x} + b = 0\} = \mathbb{R}^n - H^+ - H^- \subseteq \mathbb{R}^n$

**D. (Universal Approximation with Ridge Functions)**  
In decision theory, we strive to minimize the expected risk (defined through a loss function  $\ell$ ) of a function  $F$ .

**D. (Ridge Function)**  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a ridge function, if it can be written as  $f(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b)$  for some  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ ,  $\mathbf{w} \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ .  
**Rep.** polygonal lines by (linear function +) linear combinations of  $(-)$ - or  $(+)$ -functions

**3. Apply dimension lifting lemma to show density of the linear span of resulting ridge function families  $\mathcal{G}_{++}^n$  and  $\mathcal{G}_{-+}^n$ .**

**11.5.1 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):  
 $H^+ := \{x \mid \mathbf{w}^\top \mathbf{x} + b > 0\} \subseteq \mathbb{R}^n$   
 $H^- := \{x \mid \mathbf{w}^\top \mathbf{x} + b < 0\} \subseteq \mathbb{R}^n$   
 $H^0 := \{x \mid \mathbf{w}^\top \mathbf{x} + b = 0\} = \mathbb{R}^n - H^+ - H^- \subseteq \mathbb{R}^n$

**D. (Loss Function)**  
In decision theory, we strive to minimize the expected risk (defined through a loss function  $\ell$ ) of a function  $F$ .

**D. (Expected Risk of a Function F)**  
 $F^* = \arg \min_F \sum_{y \in \mathcal{Y}} \int_{\mathcal{X}} \ell(y, F(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x}$   
 $= \arg \min_F \mathbb{E}_{\mathbf{x} \sim \mathcal{X}} [\ell(Y, F(\mathbf{x}))]$   
 $\mathcal{R}^*(F)$  expected risk of  $F$

**11.5**

### 13.3 Convolution via Matrices

Represent the input signal, the kernel and the output as *vectors*. Copy the kernel as columns into the matrix offsetting it by one more time (gives a band matrix (special case of Toeplitz matrix)). Then the convolution is just a matrix-vector product.

### 13.4 Why to use Convolutions in DL

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

Many signals obey translation invariance, so we'd like to have translation invariant feature maps. If the relationship of translation invariance is given in the input-output relation then this is perfect.

### 13.5 Border Handling

There are different options to do this

**D. (Padding of  $p$ )** Means we extend the image (or each dimension) by  $p$  on both sides (so  $+2p$ ) and just fill in a constant there (e.g., zero).

**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

### 13.6 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  $W^l$  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.

This is what gave rise to the inception module:

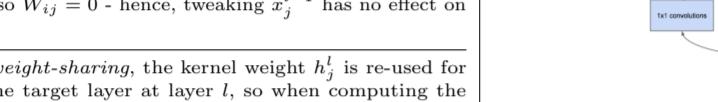
### D. (Dimension Reduction)

$m$  channels of a  $1 \times 1 \times k$  convolution  $m \leq k$ :

$$x^+ \cdot \eta = \sigma(Wx_{ij}). \quad W \in \mathbb{R}^{m \times k}.$$

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

### 13.13.5 Google Inception Network



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

- a node  $x_j^{l-1}$  may not be connected to  $x_i^l$ ,
- or a node  $x_j^{l-1}$  may be connected to  $x_i^l$  through an edge with zero weight, so  $W_{ij} = 0$  - hence, tweaking  $x_j^{l-1}$  has no effect on  $x_i^l$ .

So due to the weight-sharing, the kernel weight  $h_j^l$  is re-used for every unit in the target layer at layer  $l$ , so when computing the derivative  $\frac{\partial \mathcal{R}}{\partial h_j^l}$  we just build an additive combination of all the derivatives (note that some of them might be zero).

$$\frac{\partial \mathcal{R}}{\partial h_j^l} = \sum_{i=1}^m \frac{\partial \mathcal{R}}{\partial x_i^l} \frac{\partial x_i^l}{\partial h_j^l}$$

### Backpropagations of Convolutions as Convolutions

$y^{(l)}$  output of  $l$ -th layer  $y^{(l-1)}$  output of  $(l-1)$ -th layer / input to  $l$ -th layer w convolution filter  $\frac{\partial \mathcal{R}}{\partial y^{(l)}}$  known  $y^{(l+1)} = y^{(l)} * w$

$$\begin{aligned} \frac{\partial \mathcal{R}}{\partial w_i} &= \sum_k \frac{\partial \mathcal{R}}{\partial y_k^{(l)}} \frac{\partial y_k^{(l)}}{\partial w_i} \\ &= \sum_k \frac{\partial \mathcal{R}}{\partial y_k^{(l)}} \frac{\partial}{\partial w_i} \left[ \sum_{o=p}^k y_{o-k}^{(l-1)} w_o \right] = \sum_k \frac{\partial \mathcal{R}}{\partial y_k^{(l)}} y_{k-i}^{(l-1)} \\ &= \sum_k \frac{\partial \mathcal{R}}{\partial y_k^{(l)}} y_{k-i}^{(l-1)} = \sum_k \frac{\partial \mathcal{R}}{\partial y_k^{(l)}} \text{rot180}(y^{(l-1)})_{k-i} \\ &= \left( \frac{\partial \mathcal{R}}{\partial y^{(l)}} * \text{rot180}(y^{(l-1)}) \right)_i \end{aligned}$$

The derivative  $\frac{\partial \mathcal{R}}{\partial y^{(l)}}$  is analogous.

Note that we just used generalized indices  $i, k, o$  which may be multi-dimensional.

This example omits activation functions and biases, but that could be easily included with the chain rule.

### D. (Rotation180) $y_i := \text{rot180}(x_i) = x_{-(i)}$

### 13.7 Efficient Comp. of Convolutional Activities

A naive way to compute the convolution of a signal of length  $n$  and a kernel of length  $m$  gives an effort of  $\mathcal{O}(mn)$ . A faster way is to transform both with the FFT and then just do element-wise multiplication (effort:  $\mathcal{O}(n \log n)$ ). However, this is rarely done in CNNs as the filters usually are small ( $m \ll n, m \approx \log(n)$ ).

### 13.8 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, further we profit from dimensionality reduction. Another thing that turns out to be useful are kernels that are learned as a low-pass filter. Hence, when we sub-sample the images most of the information is still contained.

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

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:  $(\mathcal{T}, \circ^{-1}, \text{id})$ .

### 13.10 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

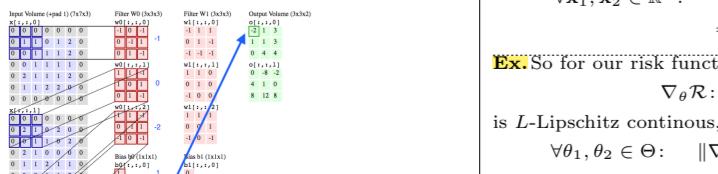
+ Increase of efficiency

### 13.11 Channels

**Ex.** Here we have

- an input signal that is 2D with 3 channels (7x7x3) (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 3x3x2 (convolution result x num convolutions)



**Ex.** We have the following chain of inclusions for functions over a closed and bounded (i.e., compact) subset of the real line.

Continuously differentiable  $\subseteq$  Lipschitz continuous  $\subseteq$  (Uniformly) continuous

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.

### 13.12 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

### 13.13 Famous CNN Architectures

#### LeNet, 1989

MNIST, 2 Convolutional Layers + 2 Fully-connected layers

#### LeNet-5

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

#### AlexNet, 2012

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

#### Inception Module

Now, a problem that arose with this ever deeper and deeper networks was 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.

#### Convex Set

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

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

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

#### Convex Function

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

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

**Com.** convex combination of two points  $\leq$  evaluation of convex combination of two points.

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

#### Every local optimum of a convex function is a global optimum

#### Operations that Preserve Convexity

- $f$  is convex if and only if  $f$  is convex
- nonnegative weighted sums
- point-elementwise maximum  $\max(f_1(x), \dots, f_n(x))$
- composition with non-decreasing function, e.g.  $e^{f(x)}$
- composition with affine mapping:  $f(Ax + b)$
- restriction to a line of convex set domain

#### Strictly Convex Function

$f$  is called *strictly convex* if

$$\forall x, x' \in S, x \neq x' \lambda \in [0, 1]: \quad f(\lambda x + (1 - \lambda)x') < \lambda f(x) + (1 - \lambda)f(x')$$

**Com.** The concept of strict convexity extends and parametrizes the notion of strict convexity. A strongly convex function is also strictly convex, but not vice versa. Notice how the definition of strong convexity approaches the definition for strict convexity as  $\mu \rightarrow \infty$ , and is related to the definition of a convex function when  $\mu = 0$ . Despite this, functions exist that are strictly convex, but are not strongly convex for any  $\mu > 0$ .

So it turns out that the non-convexity is actually not so much of an issue. It turns out that when we go to very high dimensions, the number of local minima VS the number of saddle points (gradient is zero, but non-optimal) is very small - so we're much more likely to end up in a saddle point. However, in practice, if we do SGD there is some stochasticity that will make our gradient move. Then, after waiting for a while we'll exit the saddle point.

Now, next we'll look at the insights that were gained by the paper of Saxe.

#### Optimization Challenges in NNs: Local Minima

At the beginning people were happy when they were doing convex optimization because there was a single optimum and it was reached quickly. And then when people started using non-convex optimization they were afraid of getting into non-optimal local minima and getting stuck there.

Neural network cost functions can have many local minima and/or saddle points - and this is typical. Gradient descent can get stuck. Questions that have been looked at are

- Are local minima a practical issue? Sometimes not: Gori & Tesi, 1992
- Do local minima even exist? Sometimes not (auto encoder): Baluja & Hornik, 1989
- Are local minima typically worse? often not (large networks): Arora et al., 2015
- Can we understand the learning dynamics? Deep linear case has similarities with non-linear case, e.g., Saxe et al. 2013

#### Strongly Convex Function

A differentiable function  $f$  is called *strongly convex* if

$$\forall x, x' \in S, x \neq x' \lambda \in [0, 1]: \quad f(\lambda x + (1 - \lambda)x') < \lambda f(x) + (1 - \lambda)f(x')$$

**Com.** Neural network cost functions are *strongly convex*.

So we can express the risk as

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

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

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

and the remaining terms, will be negative (as defined by the Taylor sum)

So a typical remedy for first-order methods is to take very small step sizes  $\eta$ .

However, things become even stranger because of the curvature. As we can see, the gradient norm gets larger and larger as we move along the curve. And as we can see, starting 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!

So we have that

$$\mathcal{R}(\theta(t)) - \mathcal{R}^* \leq \frac{2L}{t+1} \|\theta(t) - \theta^*\|^2 \in \mathcal{O}(t^{-1})$$

and

the remaining terms, will be negative (as defined by the Taylor sum)

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

The isometric balls illustrate the regularization loss (L2) for any choice of  $\theta$  (or  $w$ ), and the ellipsoid curves illustrate the risk (for a parabolic risk). So  $\tilde{w}$  is the point with the least loss for its specific regularization loss. As we can see, at that point

- downwards 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 not reduced that much.
- from right to left (starting at  $w^*$ ) the risk has a very low eigenvalue, and hence  $w$  is reduced much more along that dimension.

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

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

#### - 14.10.1—Regularization via Constrained Optimization

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

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

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

The simple optimization approach to this is: projected gradient descent

$$\theta(t+1) = \Pi_r(\theta(t) - \eta \nabla \mathcal{R}), \quad \Pi_r(\mathbf{v}) := \min \left\{ \frac{r}{\|\mathbf{v}\|}, \mathbf{v} \right\}$$

So we're essentially clipping the weights.

Actually, for each  $\lambda^l$  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 (so use row-norms for the weight matrices).

This had some practical success in stabilizing the optimization.

#### - 14.10.2—Early Stopping

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

Since the weights are initialized to small values (and grow and grow to fit/overfit) we're kindof clipping/constraining the weight sizes by stopping the learning process earlier.

Let's analyze the situation closer: If we study the gradient descent trajectories through a quadratic approximation of the loss around the optimal set of parameters  $\theta^*$ . We've derived previously already (and show it here again with slightly different notation) that:

$$\nabla \mathcal{R}[\theta_0] \approx \nabla \mathcal{R}[\theta_0] + \mathbf{J} \nabla \mathcal{R}[\theta_0] (\theta_0 - \theta^*) = \mathbf{H}(\theta_0 - \theta^*).$$

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

So (as seen previously) we have that

$$\theta(t+1) = \theta(t) - \eta \nabla \mathcal{R}[\theta(t)] \approx \theta(t) - \eta \mathbf{H}(\theta(t) - \theta^*).$$

Now, subtracting  $\theta^*$  on both sides gives us

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

Now we'll use the same trick as before that we can diagonalize the hessian  $\mathbf{H}$  as it's s.p.s.d., so  $\mathbf{H} = \mathbf{Q} \mathbf{A} \mathbf{Q}^\top$ . Inserting this gives us:

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

Now let's have a look at everything w.r.t. the eigenbasis of  $\mathbf{H}$ , let's define  $\tilde{\theta} = \mathbf{Q}^\top \theta$ . Then

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

Now, assuming  $\theta(0) = \mathbf{0}$  (and inserting and using it) and a small  $\eta$  ( $\forall i: |1 - \eta \lambda_i| < 1$ ) one gets explicitly

$$\tilde{\theta}(t) = \tilde{\theta}^* - (\mathbf{I} - \eta \mathbf{A})^t \tilde{\theta}^*. \quad \rightarrow 0 \text{ with upper ass. on eigenvalues}$$

Thus (comparing to the previous analysis) if we can choose  $t, \eta$  s.t.

$$(\mathbf{I} - \eta \mathbf{A})^t \stackrel{!}{=} \lambda(\mathbf{A} + \lambda \mathbf{I})^{-1}$$

which for  $\eta \epsilon_i \ll 1$ , and  $\epsilon_i \ll \lambda$  can be achieved approximately via performing  $t = \frac{1}{\lambda} \eta \epsilon_i$  steps.

So early stopping (up to the first order) can thus be seen as an approximate  $L_2$ -regularizer.

#### - 14.11—Dataset Augmentation

Applying some transformations to the input data such that we know that the output is not affected. E.g., for images: mirroring, slight rotations, scaling, slight shearing, brightness changes. Blows up data, but: there are approaches to incorporating this into the gradient instead of the input data.

#### - 14.11.1—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  $\mathbf{x}$  and apply the transformation  $\mathbf{x}' := f(\mathbf{x})$ . Then for our neural network  $F$ :

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

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

E.g., NNs where the first layer is a convolution are invariant to image translation. Hence, it would make no sense to augment the dataset of images with translations. It also saves computation and memory not to do this. So if we have an architecture that is invariant to certain dataset augmentations the augmentations become obsolete. So, if you can, choose an invariant architecture to make your life easier in the first place.

#### - 14.11.2—Injection of Noise

If we have a lot of data, but only a few datapoints are labeled. Then semi-supervised training may become useful. You may build a generative model or an autoencoder to learn how to represent your data (learn features). Then, we train a supervised model on top of these representations.

#### - 14.11.4—Multi-Task Learning

If we have different tasks that we may want to solve, we may share the intermediate representations across the tasks and then learn jointly (i.e., minimize the combined objective). A typical architecture would be to share the low-level representations, learn the high-level representations per task.

#### - 14.12—Dropout

**Dropout idea:** randomly “drop” subsets of the units in the network.

So more precisely, we'll define a “keep” probability  $\pi_i^t$  for unit  $i$  in layer  $t$ .

- typically:  $\pi_i^0 = 0.8$  (inputs),  $\pi^{T \geq 1} = 0.5$  (hidden units)
- realization: sampling bit mask and zeroing out activations
- effectively defines an exponential ensemble of networks (each of which is a sub-network of the original one), just that we sample these models at training-time (instead of during prediction) and we share the parameters
- all models share the same weights
- standard backpropagation applies.

This prevents complex co-adaptation, in which a feature detector is also helpful for producing the correct answer given the combinatorially large variety of internal contexts in which it must operate. (Hinton et al., 2012). This enforces the features to be redundant (not too specific about one thing in the image) and also to build on top of all the features of the previous layer (since we never know if some are absent).

**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 forced to start).

Equivalent to: adding multiplicative noise to weights or training exponentially many sub-networks  $\sum_{i=1}^n \binom{n}{i} = 2^n$  when  $n$  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 ensembles 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:

$$\text{Pensemble } (y | \mathbf{x}) = \sqrt{\prod_{i=1}^n P(\mu_i) P(y | \mathbf{x}_i, \mu_i)}$$

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

$$w_{ij}^t \leftarrow \pi_j^{t-1} w_{ij}^t$$

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

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

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.

#### 15. Natural Language Processing

Similarities between text and image processing: local information, Differences between text and image processing: 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

#### 15.1—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).

For this reason we'll introduce the following function

$$\text{embedding map: vocabulary } V \mapsto \mathbf{r}^d \text{ (embeddings)}$$

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

word  $w \in V \rightarrow$  one-hot word  $w \in \{0, 1\}^{|V|} \rightarrow$  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  $\mathbf{R}^d$ , 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  $\mathbf{R}^d$ ).

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, . . . )

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

So early stopping (up to the first order) can thus be seen as an approximate  $L_2$ -regularizer.

#### - 14.11—Dataset Augmentation

Applying some transformations to the input data such that we know that the output is not affected. E.g., for images: mirroring, slight rotations, scaling, slight shearing, brightness changes. Blows up data, but: there are approaches to incorporating this into the gradient instead of the input data.

#### - 14.11.1—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

$\mathbf{x}$  and apply the transformation  $\mathbf{x}' := f(\mathbf{x})$ . Then for our neural

network  $F$ :

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

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

#### - 14.11.2—Injection of Noise

If we have a lot of data, but only a few datapoints are labeled. Then semi-supervised training may become useful. You may build a generative model or an autoencoder to learn how to represent your data (learn features). Then, we train a supervised model on top of these representations.

#### - 14.11.4—Multi-Task Learning

If we have different tasks that we may want to solve, we may share the intermediate representations across the tasks and then learn jointly (i.e., minimize the combined objective). A typical architecture would be to share the low-level representations, learn the high-level representations per task.

#### - 14.12—Dropout

**Dropout idea:** randomly “drop” subsets of the units in the network.

So more precisely, we'll define a “keep” probability  $\pi_i^t$  for unit  $i$  in layer  $t$ .

- typically:  $\pi_i^0 = 0.8$  (inputs),  $\pi^{T \geq 1} = 0.5$  (hidden units)
- realization: sampling bit mask and zeroing out activations
- effectively defines an exponential ensemble of networks (each of which is a sub-network of the original one), just that we sample these models at training-time (instead of during prediction) and we share the parameters
- all models share the same weights
- standard backpropagation applies.

This prevents complex co-adaptation, in which a feature detector is also helpful for producing the correct answer given the combinatorially large variety of internal contexts in which it must operate. (Hinton et al., 2012). This enforces the features to be redundant (not too specific about one thing in the image) and also to build on top of all the features of the previous layer (since we never know if some are absent).

**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 forced to start).

So actually, what we want to do is we want to maximize the likelihood of the co-occurrences in our dataset:

$$\theta^* = \arg \max_{\theta} \prod_{(i,j) \in C_R} p_{ij}(\theta)$$

Now our approach to approximate the probability  $p_{ij}(\theta)$  as follows: it should be something that is related to the dot product of the embeddings, so

### 16.3 — Attention Mechanisms

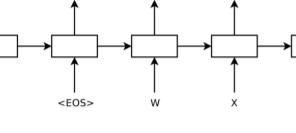
**D. (Attention Mechanisms)** offer a simple way to overcome some challenges of RNN-based memorization. With attention mechanisms we selectively attend to *inputs* or *feature representations* computed from inputs.

- RNNs: learn to encode information relevant for the future.
- Attention: selects what is relevant from the past in hindsight!

Both ideas can be combined!

If we have a sentence in English and one in German the question is how do we match one to the other. The problem with CTC was that if things are changed in order, then CTC cannot deal with it. Because the CTC doesn't process every input before it produces an output. Attention will provide a mechanism to deal with this.

So we'll see how we can do sequence to sequence learning. The idea fairly simple: Let's say we have a sequence  $ABC$  and we want to map it to  $WXYZ$ . To achieve this we'll use the so-called *encoder-decoder* architecture:



So what we'll do is

- we'll encode the sequence (e.g., sentence) into a vector, and then
- we'll decode the sequence (e.g., translate) from the vector (w/out feedback) into another sequence.

So the probability that we want to determine is

$$P(y^1, \dots, y^{T_y} | x_1, \dots, x_{T_x}, F(x^{T_x}))$$

The issue that we have here is that  $T_x$  and  $T_y$  have variable lengths, and the difference between the two lengths is not always the same. So it's very hard to match one sequence to another. Now, sequence learning will help us to function.

This expectation can be approximated by sampling.

The main problem with this is that it assumes that the two normalization constants are the same!

### 17.2 — Autoencoders

Goal: Compress the data into m-dim. ( $m \leq d$ ) representation.

**D. (Autoencoder)** any NN that aims to learn the *identity map*.

$R(\theta) = \frac{1}{2n} \sum_{i=1}^n \|x - F_\theta(x)\|_2^2 = \mathbb{E}_{x \sim p_{\text{emp}}} [\ell(x, (H \circ G)(x))]$

$$\ell(x, \hat{x}) = \frac{1}{2} \|x - \hat{x}\|_2^2$$

Typically, the network can be broken into two parts  $G$  and  $H$  such that

- $F = H \circ G \approx x \mapsto x$
- Encoder:  $G = f_1 \circ \dots \circ f_n: \mathbb{R}^m \rightarrow \mathbb{R}^m, x \mapsto z = x'$
- Decoder:  $H = f_k \circ \dots \circ f_{k+1}: \mathbb{R}^m \rightarrow \mathbb{R}^m, z \mapsto \hat{x}$
- layer  $l$  is usually the "bottleneck" layer.

Com. Just a special case of a feedforward NN, that can be trained through backpropagation.

Autoencoders provide a canonical way of *representation learning* (since NNs naturally do this). Note, how the data compression (learning compressed representation) is just a "proxy" and not the real learning objective of the network (identity function).

### 17.2.1 — Linear Autoencoding

**D. (Linear Autoencoder)**

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

$$R(\theta) = \frac{1}{2n} \sum_{i=1}^n \|x_i - DCx_i\|_2^2$$

So it's a NN with one hidden layer (no biases and linear activation functions) which will contain the compressed representation  $z = Cx \in \mathbb{R}^m$ .

### D. (Linear Autoencoder with Coupled Weights)

Then, we define  $D := C^\top$ .

### D. (Singular Value Decomposition)

Recall that the SVD of a data matrix

$$X = \begin{bmatrix} x_1 & x_2 & \cdots & x_k \end{bmatrix}$$

is of the following form:

$$X = U \operatorname{diag}(\sigma_1, \dots, \sigma_{\min(n,k)}) V^\top = \Sigma \in \mathbb{R}^{n \times k}$$

And the matrices  $U$  and  $V$  are orthogonal - so we have an orthogonal basis.

Further recall that via the SVD we can get the best rank  $k$  approximation of a linear mapping. It is also a decomposition that preserves as much of the variance (or energy) of the data for a predefined number of desired basis vectors to represent it.

### 16.4 — Recursive Networks

Good to process tree-structure, e.g., from a parser (more depth efficient  $O(\log(n))$ ). Gives a single output at the root.

$F: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

$h^n = F(h^{n-1}, h^{n-1})$

## 17 Unsupervised Learning

Here we'll look at what we can say about a distribution of  $X$ , when we have some samples  $x_1, \dots, x_N$ . Unsupervised learning is the most dangerous thing that we can do (dangerous if we don't know what we're doing). Unsupervised learning usually is hard, because we don't have a goal. The final goal of unsupervised learning is *density estimation* - so, understand the distribution that the data is coming from. Other things we might strive for is interpretability of the results we've learned about  $p(x)$ . Another key aspect of unsupervised learning is: "I don't know what I'm looking for until I find it."

### 17.1 — Density Estimation

**D. (Density Estimation)** is a standard problem in statistics and unsupervised learning. It's used to learn the distribution of the data. Classically, we use a *parametric family of densities*

$$\{p_\theta | \theta \in \Theta\}$$

to describe the set of densities that we may model. Usually, the parameters are estimated with MLE (expectation w.r.t. the empirical distribution)

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{emp}}} [\log(p_\theta(x))]$$

However, real data is rarely gaussian, laplacian, ... e.g., images. So the fact that in general we cannot solve for  $p_\theta$  for a parametric function makes this task quite complicated.

So when using a *prescribed model*  $p_\theta$  we have to

- ensure that  $p_\theta$  defines a proper density:

$$\int p_\theta(x) dx = 1$$

• and to be able to evaluate the density  $p_\theta$  at various sample points  $x$

- this may be trivial for models such as exponential families (simple formulas)
- but impractical for complex models (Markov networks, DNNs)

A typical example for an non-parametric and unnormalized model kernel-density estimation.

**D. (Kernel Density Estimator)** Let  $x_1, \dots, x_n$  be a sample, and  $k$  a kernel with bandwidth  $h > 0$  then the estimator is defined as:

$$\bar{p}_\theta(x) = \frac{1}{n} \sum_{i=1}^n k_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x - x_i}{h}\right)$$

The problem with this is that the rate of convergence is  $\log(\log(n))$  - this is extremely painfully slow. This is just a guarantee in general when we know nothing about our density.

An alternative is to use *unnormalized models* (non-parametric: the number of parameters depends on dataset size). These then represent improper density functions:

$$\bar{p}_\theta(x) = \frac{c_\theta}{\text{represented}} \cdot \frac{p_\theta(x)}{\text{unknown}} \cdot \frac{\text{normalized}}{\text{}}$$

Finding the normalization constant  $c_\theta$  might be really complicated, so we want to use relative probabilities. Further, here we cannot use the log-likelihood, because scaling up  $\bar{p}_\theta$  leads to an unbounded likelihood.

So the question still is: is there an alternative *estimation method* for unnormalized models?

What we do in practice is we do not look for the exact  $p_\theta$ , but we look for properties of  $p_\theta$ . In many cases these properties depend on our prior knowledge of  $p_\theta$ . We need to understand what the problem is in order to put the prior knowledge into the model that we want to do. This was already important in supervised learning (e.g., CNNs with several layers for images), and is even more important in unsupervised learning. We have to do the same thing there without knowing what our final goal is.

### 17.2.4 — Denoising Autoencoders

Autoencoders also allow us to separate the signal from noise: Denoising autoencoders aim to learn features of the original data representation that are robust under noise.

### D. (Score Matching (Hyvonen 2005))

$\psi_\theta := \nabla_x \log \bar{p}_\theta, \quad \psi = \nabla_x \log p$

Minimize the criterion

$$J(\theta) = \mathbb{E} [\|\psi_\theta - \psi\|^2]$$

or equivalently (by eliminating  $\psi$  by integration by parts)

$$J(\theta) = \mathbb{E} \left[ \sum_i \partial_t \psi_{t,i} - \frac{1}{2} \psi_{t,i}^2 \right]$$

This expectation can be approximated by sampling.

The main problem with this is that it assumes that the two normalization constants are the same!

### 17.2.5 — Autoencoders

Given: data points  $\{x_1, \dots, x_n\} \subset \mathbb{R}^d$

Goal: Compress the data into m-dim. ( $m \leq d$ ) representation.

**D. (Autoencoder)** any NN that aims to learn the *identity map*.

$R(\theta) = \frac{1}{2n} \sum_{i=1}^n \|x - F_\theta(x)\|_2^2 = \mathbb{E}_{x \sim p_{\text{emp}}} [\ell(x, (H \circ G)(x))]$

$$\ell(x, \hat{x}) = \frac{1}{2} \|x - \hat{x}\|_2^2$$

So we have the data  $\{x_1, \dots, x_n\} \subset \mathbb{R}^d$

Let's assume  $x_1, \dots, x_k \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{A})$ . Further let's define the data matrix  $\mathbf{X}$  as

$$\mathbf{X} = \begin{bmatrix} x_1 & x_2 & \cdots & x_k \end{bmatrix}$$

and the empirical co-variance matrix as

$$\mathbf{S} := \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{k} \mathbf{X} \mathbf{X}^\top$$

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

$$\log(P(\mathbf{X}; \mathbf{A})) = -\frac{k}{2} \operatorname{Tr}(\mathbf{SA}^{-1}) - \log(\det(\mathbf{A})) + \text{const.}$$

Note: this can be verified by using the definition of  $\mathbf{S}$ , the cyclic property of the trace, and then just write down the matrix-product as block-matrices and see what is the diagonal of the resulting matrix.

• LLE/Isomap/GPLVM (here we also try to do PCA or Factor analysis with nonlinear components (with p.w. linear components))

• Restricted Boltzmann Machine

(the idea is that  $\mathbf{Z}$  is discrete)

• Nonnegative Matrix Factorization

( $f$  "psomodel" or Bernoulli model, and both  $\mathbf{Z}$  and  $\mathbf{B}$  have to be a nonnegative matrix)

• Beta Process (aka Chinese Restaurant Process)

• Implicit Models (e.g., Generative Adversarial Networks)

(here all the information is moved to the function  $f$  instead of computing the matrices  $\mathbf{B}$  and  $\mathbf{C}$ )

### 17.4.4 — Latent Variable Models

Classically we define complex models via the *marginalization of a latent variable model*

$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\mathbf{z} | \mathbf{x}}, \Sigma_{\mathbf{z} | \mathbf{x}})$

where

$$\mu_{\mathbf{z} | \mathbf{x}} = \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} (\mathbf{x} - \mu)$$

$$\Sigma_{\mathbf{z} | \mathbf{x}} = \mathbf{I} - \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} \mathbf{W}$$

Further, if we assume that  $\Sigma = \sigma^2 \mathbf{I}$  and we let  $\sigma^2 \rightarrow 0$  (the reconstruction-error-variance for all the components is the same and we let the reconstruction error go to zero), then the following expression just reduces to the pseudo-inverses:

$$\mathbf{W}^\top (\mathbf{WW}^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{W} = \mathbf{W}^\top \in \mathbb{R}^{m \times n}$$

Consequently with the assumption of zero reconstruction error:

$$\mu_{\mathbf{z} | \mathbf{x}} \rightarrow \mathbf{W}^\top (\mathbf{x} - \mu)$$

$$\Sigma_{\mathbf{z} | \mathbf{x}} \rightarrow \mathbf{0}$$

So if we know  $\mathbf{W}$  and  $\Sigma$  is assumed to be isotropic with the error going to zero the encoding distribution gets very easy to compute.

### 17.4.3 — Dimensionality Reduction

One of the recurring things that we see in all of these models is dimensionality reduction. So we have that

$$\mathbf{X} = f(\mathbf{ZB})$$

where

- $\mathbf{X}$  is  $N \times D$ ,
- $\mathbf{Z}$  is  $N \times K$ ,
- $\mathbf{B}$  is  $K \times D$ , and
- $K \ll D$ .

So we have the data  $\mathbf{X}$  that we're trying to understand. We'll try to understand this data by a tall matrix  $\mathbf{Z}$  and a fat matrix  $\mathbf{B}$ . The tall matrix are the latent factors that we're talking about (how do we summarize the information of each sample). And the matrix  $\mathbf{B}$  is telling us how we can recover the original data from the summary. Most of the unsupervised algorithms can be captured in this general framework.