




What will it be about?

- ◆ Statistical approach to neural network learning
 - ◆ Specificity of the expectation-based learning
 - ◆ Strong law of large numbers for network learning
 - ◆ Central limit theorem for artificial neural networks
 - ◆ A central limit theorem application to network pruning
- 

Basic framework

- ◆ A MLP with n input neurons, m output neurons
- ◆ The *training pairs* $z_i = (x_i, y_i)$ with $x_i \in \mathbb{R}^n, y_i \in \mathbb{R}^m$ are viewed as realizations of *random vectors* Z_i , respectively X_i and Y_i
- ◆ All random vectors Z_i are assumed mutually independent and identically distributed with a distribution μ
 - X_i and Y_i have the marginal distributions μ_x and μ_y of μ

Assumption about moments

- ◆ Z_i and $F(X_i)$ have finite 2nd moments: $\mathbb{E}\|Z_i\|^2, \mathbb{E}\|F(X_i)\|^2 < +\infty$
 - equivalently in terms of function spaces: $Z_i \in L_2(\mu), F(X_i) \in L_2(\mu_x)$
 - for a bounded F , $F(X_i) \in L_2(\mu_x)$ already follows from $Z_i \in L_2(\mu)$
- ⇒ 1. For expectation and variance: $\mathbb{E} Z_i \in \mathbb{R}^{n+m}, \text{Var} Z_i \in \mathbb{R}^{n+m, n+m}$
- for inputs and outputs: $\mathbb{E} X_i \in \mathbb{R}^n, \mathbb{E} Y_i \in \mathbb{R}^m, \text{Var} X_i \in \mathbb{R}^{n, n}, \text{Var} Y_i \in \mathbb{R}^{m, m}$
- ⇒ 2. For conditional moments: $\mathbb{E}(Y_i|X_i) \in L_2(\mu_x), \text{Var}(Y_i|X_i) \in L_1(\mu_x)$

Expectation-based learning

- ◆ Expectation considers all possible inputs + their probability
- ◆ This learning yields weights w and biases b minimizing an expected loss $\mathbb{E}\mathcal{L}$ of network predictions $F_{(w,b)}(X)$ to outputs Y :

$$(w^*, b^*) = \arg \min_{(w,b)} \mathbb{E}_{\mu} \mathcal{L}(F_{(w,b)}(X), Y)$$

- The *most common* loss – sum of squares (SSE):

$$(w^*, b^*) = \arg \min_{(w,b)} \text{SSE} = \arg \min_{(w,b)} \mathbb{E}_{\mu} \|F_{(w,b)}(X) - Y\|^2$$

Random-sample-based learning

- ◆ Typically, $\mathbb{E}_\mu \mathcal{L}(F_{(w,b)}(X), Y)$ cannot be computed $\Leftarrow \mu$ is unknown
- ◆ But for a random sample $(x_1, y_1), \dots, (x_p, y_p)$, the mean

$\frac{1}{p} \sum_{k=1}^p \mathcal{L}(F_{(w,b)}(x_k), y_k)$ is an *unbiased estimate* of $\mathbb{E}_\mu \mathcal{L}(F_{(w,b)}(X), Y)$

- as $(x_1, y_1), \dots, (x_p, y_p)$ can serve all / some training data

- ◆ Coincides with the traditional way of learning because

$$\min_{w,b} \frac{1}{p} \sum_{k=1}^p \mathcal{L}(F_{(w,b)}(x_k), y_k) = \min_{w,b} \sum_{k=1}^p \mathcal{L}(F_{(w,b)}(x_k), y_k)$$

Specificity of SSE learning

- ◆ Notation: $\langle \cdot, \cdot \rangle_{L_2(\mu)}$ | $\|\cdot\|_{L_2(\mu)}$ – scalar product | norm in $L_2(\mu)$
- ◆
$$\text{SSE} = \left\| F_{(w,b)}(X) - Y \right\|_{L_2(\mu)}^2 = \left\| F_{(w,b)}(X) - \mathbb{E}(Y|X) \right\|_{L_2(\mu)}^2 + \left\| \mathbb{E}(Y|X) - Y \right\|_{L_2(\mu)}^2 + \left\langle F_{(w,b)}(X) - \mathbb{E}(Y|X), \mathbb{E}(Y|X) - Y \right\rangle_{L_2(\mu)}$$

$$\langle F_{(w,b)}(X) - \mathbb{E}(Y|X), \mathbb{E}(Y|X) - Y \rangle_{L_2(\mu)}$$

$$\begin{aligned} & \langle F_{(w,b)}(X) - \mathbb{E}(Y|X), \mathbb{E}(Y|X) - Y \rangle_{L_2(\mu)} = \\ & = \mathbb{E}_\mu \left(F_{(w,b)}(X) - \mathbb{E}(Y|X) \right)^\top (\mathbb{E}(Y|X) - Y) = \\ & = \mathbb{E}_{\mu_X} \left[\mathbb{E} \left(\left(F_{(w,b)}(X) - \mathbb{E}(Y|X) \right)^\top (\mathbb{E}(Y|X) - Y) \middle| X \right) \right] = \\ & = \mathbb{E}_{\mu_X} \left[\left(F_{(w,b)}(X) - \mathbb{E}(Y|X) \right)^\top \mathbb{E}(\mathbb{E}(Y|X) - Y | X) \right] = \\ & = \mathbb{E}_{\mu_X} \left[\left(F_{(w,b)}(X) - \mathbb{E}(Y|X) \right)^\top (\mathbb{E}(Y|X) - \mathbb{E}(Y|X)) \right] = 0 \end{aligned}$$

$$\|\mathbb{E}(Y|X) - Y\|_{L_2(\mu)}^2$$

$$\begin{aligned}\|\mathbb{E}(Y|X) - Y\|_{L_2(\mu)}^2 &= \mathbb{E}_\mu \|\mathbb{E}(Y|X) - Y\|^2 = \\ &= \mathbb{E}_{\mu_x} \mathbb{E}(\|\mathbb{E}(Y|X) - Y\|^2 | X) = \\ &= \mathbb{E}_{\mu_x} \left(\sum (\mathbb{E}_\mu(Y|X)_i - Y_i)^2 \mid X \right) = \\ &= \mathbb{E}_{\mu_x} \sum \text{Var}(Y|X)_{i,i} = \\ &= \mathbb{E}_{\mu_x} \text{trace Var}(Y|X)\end{aligned}$$

Specificity of SSE learning

- ◆ Notation: $\langle \cdot, \cdot \rangle_{L_2(\mu)}$ | $\|\cdot\|_{L_2(\mu)}$ – scalar product | norm in $L_2(\mu)$
- ◆ $SSE = \left\| F_{(w,b)}(X) - Y \right\|_{L_2(\mu)}^2 = \left\| F_{(w,b)}(X) - \mathbb{E}(Y|X) \right\|_{L_2(\mu)}^2 + \left\| \mathbb{E}(Y|X) - Y \right\|_{L_2(\mu)}^2 + \left\langle F_{(w,b)}(X) - \mathbb{E}(Y|X), \mathbb{E}(Y|X) - Y \right\rangle_{L_2(\mu)}$
 $= \left\| F_{(w,b)}(X) - \mathbb{E}(Y|X) \right\|_{L_2(\mu)}^2 + \mathbb{E}_{\mu_x} \text{trace Var}(Y|X)$
- ◆ Thus $\arg \min_{(w,b)} SSE = \arg \min_{(w,b)} \left\| F_{(w,b)}(X) - \mathbb{E}(Y|X) \right\|_{L_2(\mu)}^2$, and if exist $(w,b)_{(Y|X)}$ such that $F_{(w,b)_{(Y|X)}}(X) = \mathbb{E}(Y|X)$, then $\arg \min_{(w,b)} SSE = (w,b)_{(Y|X)}$

Random-sample \bowtie expectation

- ◆ Expectation-based learning is not common because the distribution of learning samples is typically unknown
- ◆ But random sample empirical *mean estimates expectation*
 1. in an *unbiased* way: $\mathbb{E}_{\mu} \frac{1}{p} \sum_{k=1}^p \mathcal{L}(F(x_k), y_k) = \mathbb{E}_{\mu} \mathcal{L}(F(X), Y)$
 2. in a *consistent* way: $\frac{1}{p} \sum_{k=1}^p \mathcal{L}(F(x_k), y_k) \rightarrow \mathbb{E}_{\mu} \mathcal{L}(F(X), Y)$

Laws of large numbers

◆ The consistence property, that $\frac{1}{p} \sum_{k=1}^p \mathcal{L}(F(x_k), y_k)$ converges to $\mathbb{E}_\mu \mathcal{L}(F(X), Y)$ is called law of large numbers.

◆ *Weak* law: convergence of random variables in *probability*

$$\forall \varepsilon > 0: \lim_{p \rightarrow \infty} \mu \left(\left| \frac{1}{p} \sum_{k=1}^p \mathcal{L}(F(x_k), y_k) - \mathbb{E}_\mu \mathcal{L}(F(X), Y) \right| > \varepsilon \right) = 0$$

◆ *Strong* law (\implies weak law): convergence *almost everywhere*

$$\mu \left(\lim_{p \rightarrow \infty} \frac{1}{p} \sum_{k=1}^p \mathcal{L}(F(x_k), y_k) = \mathbb{E}_\mu \mathcal{L}(F(X), Y) \right) = 1$$

Complete probability space

◆ Laws of large numbers cannot be directly applied

to MLPs $\Leftarrow \sum_{k=1}^p \mathcal{L}(F_{(w,b)}(x_k), y_k)$ changed by minimum

- therefore, for MLPs, specific *additional assumptions* are needed

◆ Let $Z_i, i \in \mathbb{N}$, be Borel-measurable, $Z_i: (\Omega, \mathcal{A}, P) \rightarrow (\mathbb{R}^{n+m}, \mathcal{B}, \mu)$

◆ The probability space (Ω, \mathcal{A}, P) is assumed being *complete*:

$$A \in \mathcal{A} \& B \subset \Omega \& (A \setminus B) \cup (B \setminus A) \subset C \in \mathcal{A} \& P(C) = 0 \implies B \in \mathcal{A}$$

Assumptions for the strong law

1. (Ω, \mathcal{A}, P) is a complete probability space
2. $(X_i, Y_i), i \in \mathbb{N}$, are *i.i.d.* (independent and identically distributed)
3. $W = \{\text{admissible } (w, b) | w - \text{weights}, b - \text{bias}\}$ is a *compact* set
4. $(\forall (w, b) \in W) \mathcal{L}(F_{(w,b)}(x), y)$ is a Borel-*measurable* function of (x, y)
5. $(\forall (x, y) \in \mathbb{R}^{n+m}) \mathcal{L}(F_{(w,b)}(x), y)$ is a *W-continuous* function of (w, b)
6. $\mathcal{L}(F_{(w,b)}(X), Y)$ has an \mathbb{R}^{n+m} -integrable *majorizer* over W

Statement of the strong law

- ◆ Consider the set of expectation-based learning results

$$W^* = \left\{ (w^*, b^*) \in W \mid \mathbb{E}_\mu \mathcal{L}(F_{(w^*, b^*)}(X), Y) = \min_{(w, b)} \mathbb{E}_\mu \mathcal{L}(F_{(w, b)}(X), Y) \right\}$$

+ random-sample-based learning results for $(x_i, y_i)_{p=1}^\infty$

$$(\forall p \in \mathbb{N}) (\hat{w}_p, \hat{b}_p) = \arg \min_{w, b} \sum_{k=1}^p \mathcal{L}(F_{(w, b)}(x_k), y_k)$$

- ◆ Then $(\hat{w}_p, \hat{b}_p)_{p=1}^\infty$ converges *almost everywhere to W^**

$$\mu \left(\liminf_{p \rightarrow \infty} \inf_{(w^*, b^*) \in W^*} \|(\hat{w}_p, \hat{b}_p) - (w^*, b^*)\| = 0 \right) = 1$$

Network pruning

- ◆ *Removing connections* from fully connected networks
 - decreases the risk of overtraining + computational costs
- ◆ If all input connections / all output connections of a *hidden neuron* h are pruned, then h is removed
- ◆ *Formalised*: $S(w, b) = 0$ with a 0/1-valued matrix S , rows contain for the 1 connection / for neuron's all connections + bias

Statistical approach to pruning

- ◆ Because (\hat{w}_p, \hat{b}_p) that results from learning is only an estimate (unbiased + consistent) of (w^*, b^*) , what we actually need is to know *whether* $S(w^*, b^*) = 0$
 - cannot be directly checked $\Leftarrow (w^*, b^*)$ is not known
- ◆ Statistical approach to checking statements for estimated values:
hypotheses testing using their *estimator* $((\hat{w}_p, \hat{b}_p))$

Hypotheses testing recalled

- ◆ Testing a null hypotheses H_0 against H_1 : checking whether $T \in \mathfrak{C}$
 - T - *test statistics*: random variable with some assumed distribution
 - \mathfrak{C} - *critical set*: $\mathfrak{C} \subset \text{Val} T$ with $H_0 \Rightarrow P(T \in \mathfrak{C} | H_0 \vee H_1) \leq \alpha$ - significance
- ◆ The assumed T distribution can always asymptotically rely on the normality of $\frac{1}{\sqrt{p}} \sum_{k=1}^p \mathcal{L}(F(x_k), y_k) \Leftarrow CLT$ (central limit theorem)
 - not directly applicable $\Leftarrow \sum_{k=1}^p \mathcal{L}(F_{(wb)}(x_k), y_k)$ changed by minimum

CLT for MLPs: assumptions 1.– 6.

1. (Ω, \mathcal{A}, P) is a complete probability space
2. $(X_i, Y_i), i \in \mathbb{N}$, are *i.i.d.* (independent and identically distributed)
3. $W = \{\text{admissible } (w, b) | w - \text{weights}, b - \text{bias}\}$ is a *compact* set
4. $W^* = \{(w^*, b^*)\}$ with (w^*, b^*) an *inner* point of W
5. $(\forall (w, b) \in W) \mathcal{L}(F_{(w,b)}(x), y)$ is a Borel-*measurable* function of (x, y)
6. $\mathcal{LZ}(F_{(w,b)}(X), Y)$ has an \mathbb{R}^{n+m} -integrable *majorizer* over W

CLT for MLPs: auxiliary notation

- ◆ $\nabla_{(w,b)} \mathcal{L} = \nabla_{(w,b)} \mathcal{L}(F_{(w,b)}(X), Y)$: a random vector such that $(\forall (x, y) \in \mathbb{R}^{n+m}) \nabla_{(w,b)} \mathcal{L} =$ the *gradient* of $\mathcal{L}(F_{(w,b)}(x), y)$ w.r. to (w, b)
- ◆ $\nabla_{(w,b)}^2 \mathcal{L}(F_{(w,b)}(X), Y)$: a random matrix such that $(\forall (x, y) \in \mathbb{R}^{n+m}) \nabla_{(w,b)}^2 \mathcal{L}(F_{(w,b)}(x), y) =$ the *Hessian* of $\mathcal{L}(F_{(w,b)}(x), y)$ w.r. to (w, b)

CLT for MLPs: assumptions 7.- 12.

7. $(\forall x, y) \mathcal{L}(F_{(w,b)}(x), y)$ has W -continuous Hessian w.r. to (w, b)
8. The matrix A^* defined $A^* = \mathbb{E}_{\mu} \left(\nabla_{(w^*, b^*)}^2 \mathcal{L}(F_{(w,b)}(X), Y) \right)$ is regular
9. $\nabla_{(w,b)}^2 \mathcal{L}(F_{(w,b)}(x), y)$ has an \mathbb{R}^{n+m} -integrable majorizer over W
10. The matrix B^* defined $B^* = \mathbb{E}_{\mu} \left(\nabla_{(w,b)} \mathcal{L}^\top \nabla_{(w,b)} \mathcal{L} \right)$ is regular
11. $\|\mathcal{L}(F_{(w,b)}(x), y)\|^2$ has an \mathbb{R}^{n+m} -integrable majorizer over W
12. A $\{0,1\}$ -valued matrix S has $s = \text{rank } S$ rows

CLT for MLPs: auxiliary notation

- ◆ $\nabla_{(w,b)} \mathcal{L} = \nabla_{(w,b)} \mathcal{L}(F_{(w,b)}(X), Y)$: a random vector such that $(\forall(x,y) \in \mathbb{R}^{n+m}) \nabla_{(w,b)} \mathcal{L} =$ the *gradient* of $\mathcal{L}(F_{(w,b)}(x), y)$ w.r. to (w, b)
- ◆ $\nabla_{(w,b)}^2 \mathcal{L}(F_{(w,b)}(X), Y)$: a random matrix such that $(\forall(x,y) \in \mathbb{R}^{n+m}) \nabla_{(w,b)}^2 \mathcal{L}(F_{(w,b)}(x), y) =$ the *Hessian* of $\mathcal{L}(F_{(w,b)}(x), y)$ w.r. to (w, b)
- ◆ $\hat{A}_p = \frac{1}{p} \sum_{i=1}^p \nabla_{(w,b)}^2 \mathcal{L} \left(F_{(\hat{w}_p, \hat{b}_p)}(x_i), y_i \right)$
- ◆ $\hat{B}_p = \frac{1}{p} \sum_{i=1}^p \nabla_{(w,b)} \mathcal{L} \left(F_{(\hat{w}_p, \hat{b}_p)}(x_i), y_i \right) \mathcal{L} \left(F_{(\hat{w}_p, \hat{b}_p)}(x_i), y_i \right)^\top$

CLT for MLPs: exact covariance

◆ $\left(\sqrt{p} \left((\hat{w}_p, \hat{b}_p) - (w^*, b^*) \right) \right)_{p=1}^{\infty}$ converges to the *distribution* $N(0, C^*)$,

the covariance matrix of which is $C^* = A^{*-1} B^* A^{*-1}$

◆ If $S(w^*, b^*) = 0$, then $\left(\sqrt{p} S(\hat{w}_p, \hat{b}_p) \right)_{p=1}^{\infty}$ converges to $N(0, SC^*S^\top)$

◆ If $S(w^*, b^*) = 0$, then the quadratic forms of $\left(\sqrt{p} S(\hat{w}_p, \hat{b}_p) \right)_{p=1}^{\infty}$

$$\left(p(\hat{w}_p, \hat{b}_p)^\top S^\top (SC^*S^\top)^{-1} S(\hat{w}_p, \hat{b}_p) \right)_{p=1}^{\infty}$$

converge to the *distribution* χ_s^2 with s degrees of freedom

CLT for MLPs: estimated covariance

- ◆ Define an estimate $\hat{C}_p = \begin{cases} \hat{A}_p^{-1} \hat{B}_p \hat{A}_p^{-1} & \text{if } p \in \mathbb{N}, \hat{A}_p \text{ is regular} \\ \hat{B}_p & \text{if } p \in \mathbb{N}, \hat{A}_p \text{ is singular} \end{cases}$
- ◆ Then $\hat{C}_p \rightarrow C^*$ in probability $\left(\Rightarrow \lim_{p \rightarrow \infty} \mu(\hat{C}_p \text{ is regular}) = 1 \right)$
- ◆ If $S(w^*, b^*) = 0$, then also the quadratic forms $\left(p(\hat{w}_p, \hat{b}_p)^\top S^\top (S \hat{C}_p S^\top)^{-1} S(\hat{w}_p, \hat{b}_p) \right)_{p=1}^{\infty}$ converge to the *distribution* χ_s^2

Test procedure for $S(w^*, b^*) = 0$

1. For given observations $(x_1, y_1), \dots, (x_p, y_p)$, get [= compute] (\hat{w}_p, \hat{b}_p)
2. For $i = 1, \dots, p$, get $\nabla_{(w,b)} \mathcal{L}(F_{(\hat{w}_p, \hat{b}_p)}(x_i), y_i)$ and $\nabla_{(w,b)}^2 \mathcal{L}(F_{(\hat{w}_p, \hat{b}_p)}(x_i), y_i)$
3. Get \hat{A}_p, \hat{B}_p and check whether any is singular
4. Then the test cannot proceed, else get \hat{C}_p
5. Get $p(\hat{w}_p, \hat{b}_p)^\top S^\top (S \hat{C}_p S^\top)^{-1} S(\hat{w}_p, \hat{b}_p)$ and compare with the distribution χ_S^2
6. If $\downarrow >$ the quantile $\chi_S^2(1 - \alpha), \alpha \in (0,1)$, *reject* $S(w^*, b^*) = 0$