

Statistical Data Analysis

Dr. Jana de Wiljes

1. Dezember 2021

Universität Potsdam

III-posedness and Regularization

If the least squares problem is ill-posed, i.e., solution does not exist or is unstable.

Small perturbations in \mathbf{y} or \mathbf{X} yield large perturbations in β

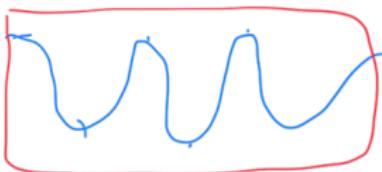
$$\tau_u(\mathbf{x}) = \rho + 1$$

Solve regularized problem: For some $\lambda > 0$ and matrix \mathbf{G}

$$\min_{\beta} \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}^T\|^2 + \frac{\lambda}{2} \|\mathbf{G}\beta\|^2$$

$$A\mathbf{x} = \mathbf{b} + \varepsilon(x)$$

$\text{rank}(A) = 2$, 3 unknown



Iterative Methods

Iterative Solvers for Least-Squares Regression

So far: Given $\mathbf{y} \in \mathbb{R}^n$, solve

$$\min_{\beta} \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|$$

directly using $\beta^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$. Here

$$\mathbf{X} \in \mathbb{R}^{n \times (p+1)} \quad \text{and} \quad \beta \in \mathbb{R}^{(p+1)}$$

Problems:

1. Generating $\mathbf{X}^\top \mathbf{X}$ and solving normal equations is too costly for large-scale problems.
2. Exact solution not useful when problem is ill-posed \rightsquigarrow add explicit regularization or do so implicitly by early stopping.

Iterative methods that avoid working with $\mathbf{X}^\top \mathbf{X}$

- Steepest descent
- Conjugate gradient for least-squares (CGLS)

Excellent references: Numerical Optimization [4], iterative linear algebra [5], general introduction [1]



Iterative Methods

General idea - obtain a sequence $\beta_1, \dots, \beta_j, \dots$ that converges to least-squares solution β^*

$$\beta_j \longrightarrow \beta^*, \quad \text{for } j \rightarrow \infty.$$

How fast does the sequence converge? Assume

$$\|\beta_{j+1} - \beta^*\| < \gamma_j \|\beta_j - \beta^*\|$$

where all $\gamma_j < 1$. Then

- If γ_j is bounded away from 0 and 1 the convergence is linear
- If $\gamma_j \rightarrow 0$ the convergence is superlinear
- If $\gamma_j \rightarrow 1$ the convergence is sublinear

The sequence converges quadratically if γ_j is bounded away from 0 and 1 and

$$\|\beta_{j+1} - \beta^*\| < \gamma_j \|\beta_j - \beta^*\|^2$$

Steepest Descent for Least-Squares

Consider now

$$\phi(\beta) = \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|^2$$

with

$$\nabla_{\beta} \phi(\beta) = \mathbf{X}^T (\mathbf{X}\beta - \mathbf{y})$$

$$\mathbf{X}^T \mathbf{X}$$

Steepest descent direction is $\mathbf{d}_j = (\mathbf{X}^T \mathbf{y} - \mathbf{X}\beta_j)$ and

$$\beta_{j+1} = \beta_j + \alpha_j \mathbf{d}_j$$

How to choose α_j ?

Idea: Minimize ϕ along direction \mathbf{d}_j

$$\alpha_j = \underset{\alpha}{\operatorname{argmin}} \phi(\beta_j + \alpha \mathbf{d}_j) = \underset{\alpha}{\operatorname{argmin}} \frac{1}{2} \|\alpha \mathbf{X} \mathbf{d}_j - \mathbf{r}_j\|^2$$

with residual $\mathbf{r}_j = \mathbf{y} - \mathbf{X}\beta_j$.

$$\mathbf{X}\beta_0$$

$$\mathbf{X}^T \mathbf{X}$$

$$\beta_j$$

$$\boxed{\mathbf{X}^T \mathbf{X}}$$

This leads to simple quadratic equation in 1D whose solution is

$$\alpha_j = \frac{\mathbf{r}_j^T \mathbf{X} \mathbf{d}_j}{\|\mathbf{X} \mathbf{d}_j\|^2}$$

Algorithm: Steepest Descent for Least-Squares

for $j = 1, \dots$

till?

β_1 = random

- Compute residual $r_j = y - X\beta_j$
- Compute the SD direction $d_j = X^\top r_j$
- Compute step size $\alpha_j = \frac{r_j^\top X d_j}{\|X d_j\|^2}$
- Take the step $\beta_{j+1} = \beta_j + \alpha_j d_j$

Converges linearly, i.e.,

$$\|\beta_{j+1} - \beta^*\| < \gamma \|\beta_j - \beta^*\|$$

with

$$\gamma \approx \left| \frac{\kappa - 1}{\kappa + 1} \right|$$

Here, κ depends on condition number of X , i.e.,

$$\kappa \approx \frac{\sigma_{\max}^2}{\sigma_{\min}^2}$$

Can be painfully slow for ill-conditioned problems

$$\hat{\beta} = \beta_{100}$$

Accelerating Steepest Descent: Post-Conditioning

Idea: Improve convergence by transforming the problem

$$\phi(\beta) = \frac{1}{2} \| \mathbf{X} \mathbf{S} \mathbf{S}^{-1} \beta - \mathbf{y} \|^2$$

Here: \mathbf{S} is invertible

Solve in two steps:

1. Set $\mathbf{z} = \mathbf{S}^{-1} \beta$ and compute

$$\mathbf{z}^* = \underset{\mathbf{z}}{\operatorname{argmin}} \frac{1}{2} \| \mathbf{X} \mathbf{S} \mathbf{z} - \mathbf{y} \|^2$$

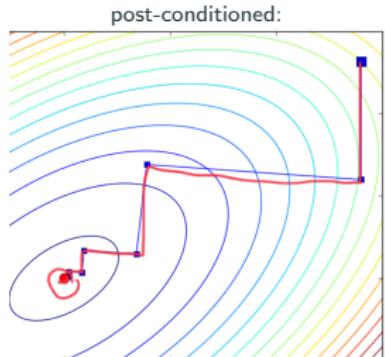
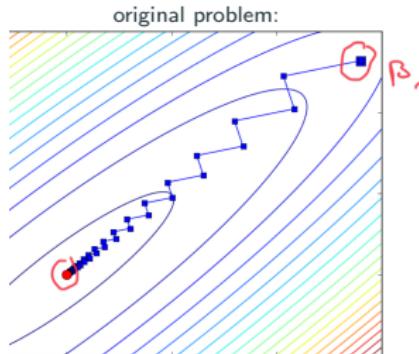
2. Then $\beta = \mathbf{S} \mathbf{z}$.

Pick \mathbf{S} such that $\mathbf{X} \mathbf{S}$ is better conditioned.

$$\boxed{\mathbf{X}^\top \mathbf{X}}$$

$$\mathbf{A} \mathbf{x} \approx \mathbf{b}$$

$$\boxed{Q}$$



Conjugate Gradient Method for Least-Squares

CG is designed to solve quadratic optimization problems

$$\min_{\beta} \frac{1}{2} \beta^T H \beta - b^T \beta$$

with H symmetric positive definite. In our case

$$\operatorname{argmin}_{\beta} \frac{1}{2} \|X\beta - y\|^2 = \operatorname{argmin}_{\beta} \frac{1}{2} \beta^T \underbrace{X^T X}_{=H} \beta - \underbrace{y^T X \beta}_{=b^T}$$

CG improves over SD by using previous step (not a memory-less method) and constructing a basis for the solution.

Facts:

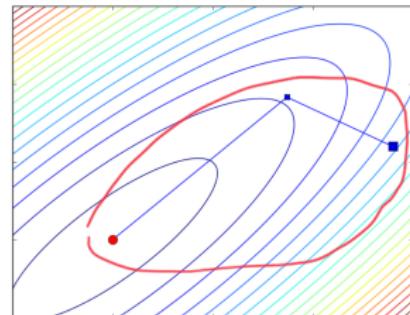
- terminates after at most n steps (in exact arithmetic)
- good solutions for $j \ll n$
- convergence $\gamma_j \approx \left| \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right|^j$

Conjugate Gradient Least-Squares

- Uses the structure of the problem to obtain stable implementation
- Typically converges much faster than SD
- Accelerate using post conditioning

$$\min_{\beta} \frac{1}{2} \|\mathbf{X}\mathbf{S}\mathbf{S}^{-1}\beta - \mathbf{y}\|^2$$

- Faster convergence when eigenvalues of $\mathbf{S}^T \mathbf{X}^T \mathbf{X} \mathbf{S}$ are clustered.



Iterative Regularization

Consider

$$\min_{\beta} \|\mathbf{X}\beta - \mathbf{b}\|^2$$

- Assume that \mathbf{X} has non-trivial null space
- The matrix $\mathbf{X}^\top \mathbf{X}$ is not invertible
- Can we still use iterative methods (CG, CGLS, ...)?

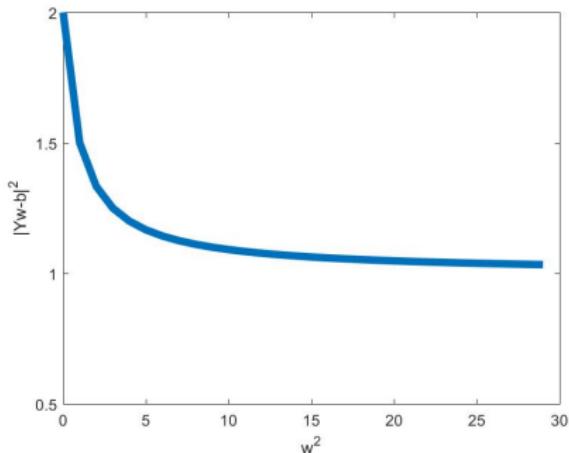
What are the properties of the iterates?

Excellent introduction to computational inverse problems [2, 6, 3]

Iterative Regularization: L-Curve

The CGLS algorithm has the following properties

- For each iteration $\|\mathbf{X}\beta_k - \mathbf{y}\|^2 \leq \|\mathbf{X}\beta_{k-1} - \mathbf{y}\|^2$
- If starting from $\beta = 0$ then $\|\beta_k\|^2 \geq \|\beta_{k-1}\|^2$
- β_1, β_2, \dots converges to the minimum norm solution of the problem
- Plotting $\|\beta_k\|^2$ vs $\|\mathbf{X}\beta_k - \mathbf{y}\|^2$ typically has the shape of an L-curve



Cross Validation - 1

Finding good least-squares solution requires good parameter selection.

- λ when using Tikhonov regularization (weight decay)
- number of iteration (for SD and CGLS)

Suppose that we have two different "solutions"

$$\begin{aligned}\beta_1 \quad \rightarrow \quad & \|\beta_1\|^2 = \eta_1 \quad \|\mathbf{X}\beta_1 - \mathbf{y}\|^2 = \rho_1. \\ \beta_2 \quad \rightarrow \quad & \|\beta_2\|^2 = \eta_2 \quad \|\mathbf{X}\beta_2 - \mathbf{y}\|^2 = \rho_2.\end{aligned}$$

How to decide which one is better?

Cross Validation - 2

Goal: Gauge how well the model can predict new examples.

Let $\{\mathbf{x}_{\text{CV}}, \mathbf{y}_{\text{CV}}\}$ be data that is **not used** for the training

Idea: If $\|\mathbf{x}_{\text{CV}}\beta_1 - \mathbf{y}_{\text{CV}}\|^2 \leq \|\mathbf{x}_{\text{CV}}\beta_2 - \mathbf{y}_{\text{CV}}\|^2$, then β_1 is a better solution than β_2 .

When the solution depends on some hyper-parameter(s) λ , we can phrase this as bi-level optimization problem

$$\lambda^* = \operatorname{argmin}_\lambda \|\mathbf{x}_{\text{CV}}\beta(\lambda) - \mathbf{y}_{\text{CV}}\|^2,$$

$$\text{where } \beta(\lambda) = \operatorname{argmin}_\beta \frac{1}{2} \|\mathbf{x}\beta - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\beta\|^2.$$

Cross Validation - 3

To assess the final quality of the solution cross validation is not sufficient (why?).

Need a final testing set.

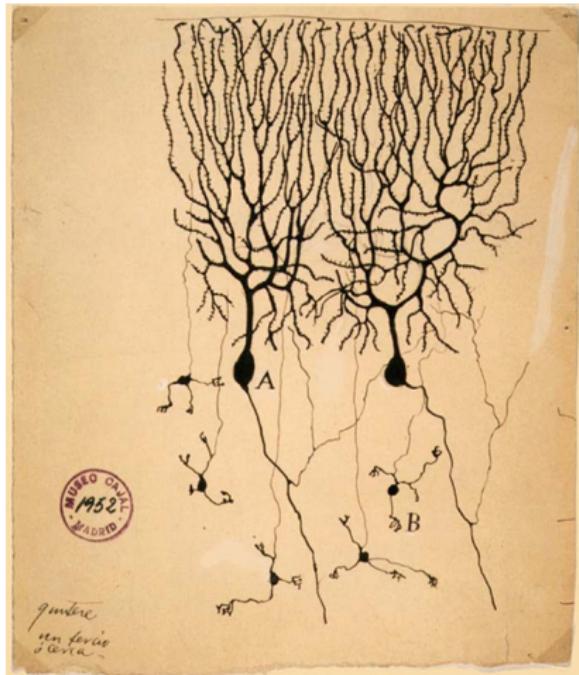
Procedure

- Divide the data into 3 groups $\{X_{\text{train}}, X_{\text{CV}}, X_{\text{test}}\}$.
- Use X_{train} to estimate $\beta(\lambda)$
- Use X_{CV} to estimate λ
- Use X_{test} to assess the quality of the solution

Important - we are not allowed to use X_{test} to tune parameters!

Neural Networks

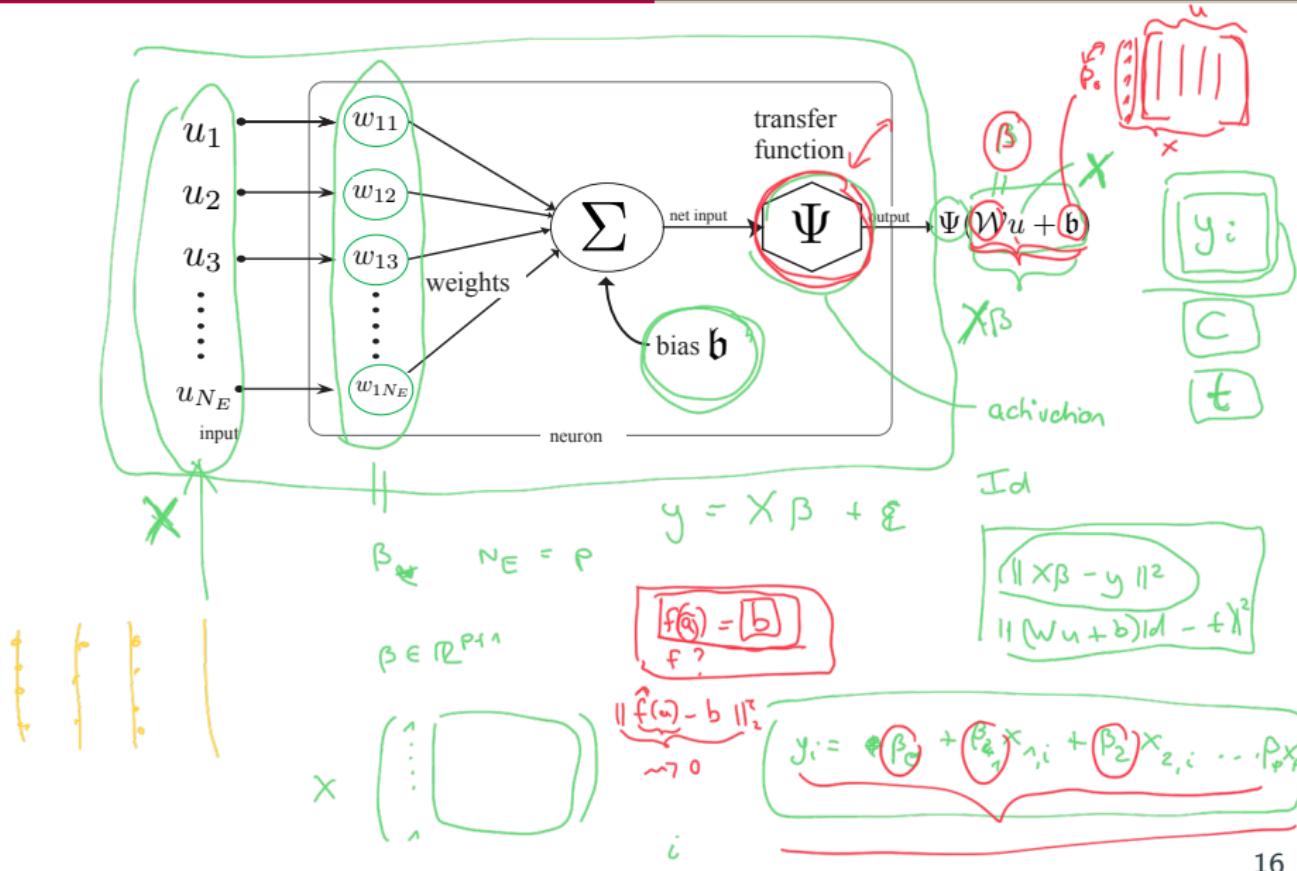
Motivation from biology



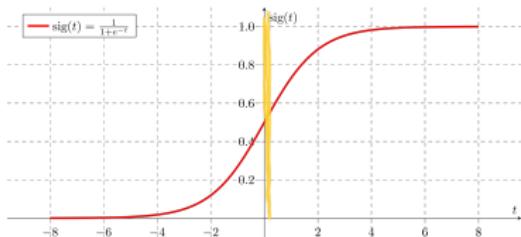
By Santiago Ramón y Cajal in 1899 see

https://de.wikipedia.org/wiki/Santiago_Ramn_y_Cajal for details

Neuron



Activation function example: sigmoid



Sigmoid function:

$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

Properties:

- Derivative:
$$\frac{1 + e^{-x} + xe^{-x}}{(1 + e^{-x})^2}$$

- $$\text{sig}'(t) = \text{sig}(t)(1 - \text{sig}(t))$$

Activation function example: ReLu



Rectified linear unit:
$$\text{ReLU}(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$$
$$= \max\{0, x\}$$

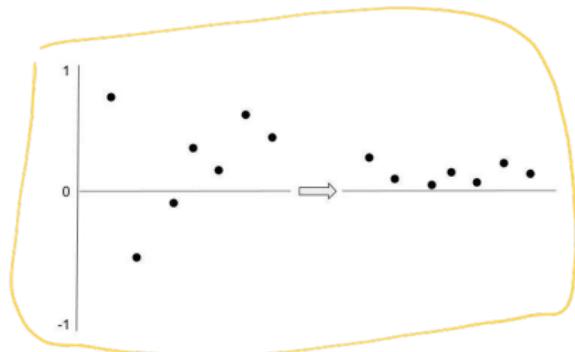
Properties:

- Derivative:

$$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \\ \text{undefined} & \text{if } x = 0 \end{cases} \quad (1)$$

- very popular for Deep RL
- Dying ReLU problem - vanishing gradient problem.

Activation function example: Softmax



Softmax:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \text{for } i = 1, \dots, K$$

and $\mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$.

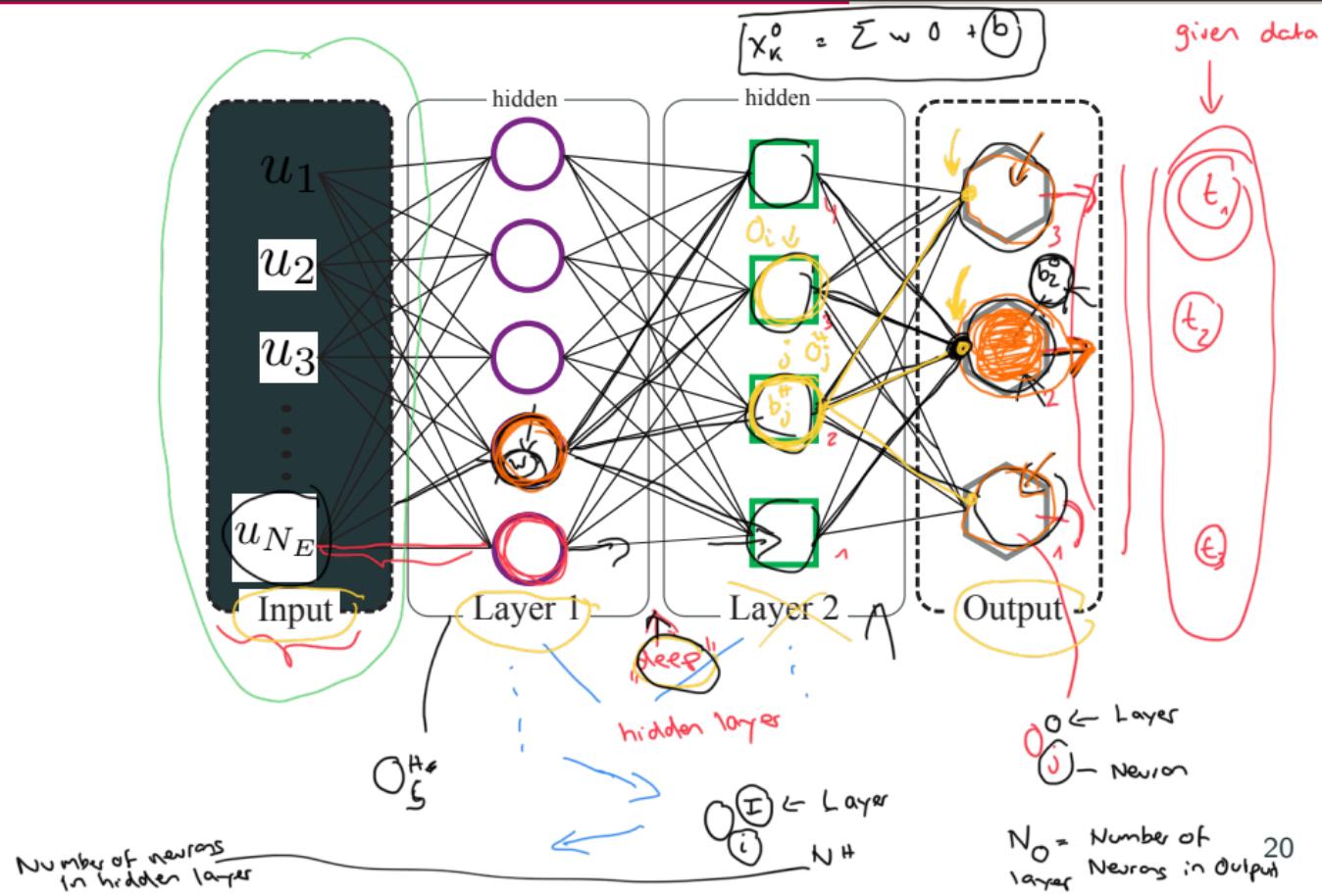
Properties:

- Derivative:

$$\left. \frac{\partial}{\partial q_k} \sigma(\mathbf{q}, i) \right| = \sigma(\mathbf{q}, i)(\delta_{ik} - \sigma(\mathbf{q}, k)). \quad (2)$$

- used in to normalize the output (map to a probability distribution)
- also used in RL to convert action values into action probabilities

Multilayer perceptron



Number of neurons
in hidden layers

N_H^1

N_H^2

N_O
 I_i ← Layer
 i ← Neuron

N_O = Number of
Neurons in Output
layer 20

Training Neural Network

1. Choose network architecture:

- activation functions - sigmoid
- hidden layers (shallow or deep) ↗
- number of neurons N_o^*, N_H, N_I
- etc.

2. Choose appropriate loss function E , e.g., least squares

3. Find minima via:

- stochastic gradient descent
- Backpropagation

The diagram illustrates a neural network layer. Inputs O_s^0 are multiplied by weights w_s and summed with bias b_s to produce the output O_s^0 . The error for each output is calculated as $(O_s^0 - t_s)^2$. The total error $E(w, b)$ is the average of these squared errors over all samples $s = 1, \dots, C$.

$$E(w, b) = \frac{1}{2} \sum_{s=1}^C (O_s^0 - t_s)^2$$

2()

$$(x^2)' = 2x$$

Stochastic Gradient Descent

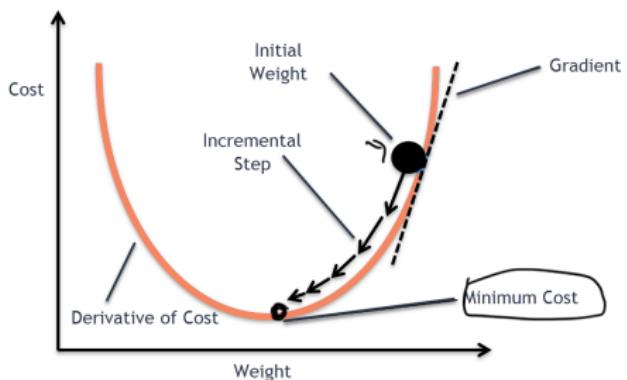
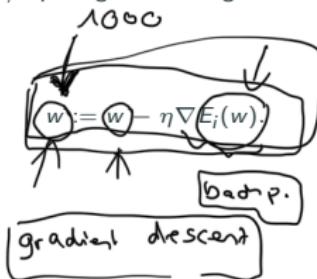


Image ref: <https://morioh.com/p/bc6bc20e9739> and
<https://medium.com/38th-street-studios/exploring-stochastic-gradient-descent-with-restarts-sgdr-fa206c38a74e>

Iterative weight improvement:



(3)

Backpropagation

$$\frac{\partial E}{\partial b_u^0} = \frac{\partial}{\partial b_u^0} \frac{1}{2} \sum_{s=1}^{N_o} (o_s^0 - t_s)^2 = (o_u^0 - t_u) \frac{\partial}{\partial b_u^0} (o_u^0 - t_u)$$

chain rule

$$= (o_u^0 - t_u) \frac{\partial}{\partial b_u^0} \text{sig}(x_u^0)$$

$$= (o_u^0 - t_u) \text{sig}(x_u^0) (1 - \text{sig}(x_u^0)) \cdot \frac{\partial x_u^0}{\partial b_u^0}$$

gradient descent step size

$$\sim b_u^0 = b_u^0 - (o_u^0 - t_u) \text{sig}(x_u^0) (1 - \text{sig}(x_u^0))$$



$$\text{Sig}(x_u^0)$$

$$\frac{\partial x_u^0}{\partial b_u^0} = \sum_{j=1}^{N_h} w_{ju}^0 o_j^0 + b_u^0$$

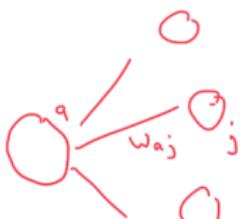
Backpropagation

$$\frac{\partial E}{\partial b_j^H} = \frac{\partial}{\partial b_j^H} \frac{1}{2} \sum_{s=1}^{N_o} (o_s^o - t_s)^2$$

$$= \sum_{s=1}^{N_o} (o_s^o - t_s) \frac{\partial}{\partial b_j^H} (o_s^o - t_s)$$

$\text{sig}(x_s^o)$

$$= \sum_{s=1}^{N_o} (o_s^o - t_s) \text{sig}(x_s^o) (1 - \text{sig}(x_s^o)) \frac{\partial x_s^o}{\partial b_j^H}$$



$a \rightarrow j$

$$\frac{\partial E}{\partial w_{ij}^H}$$

$$\sum_{a=1}^{N_I} o_a^I w_{aj}^H + b_j^H$$

$$\sum_{i=1}^{N_H} w_{is}^H o_i^H + b_s^o$$

j

$$\text{sig}(x_j)$$

$$\sum_{a=1}^{N_I} o_a^I w_{aj}^H$$

$$\frac{\partial x_s^o}{\partial o_j^H}$$

$$\frac{\partial o_j^H}{\partial w_{ij}^H}$$

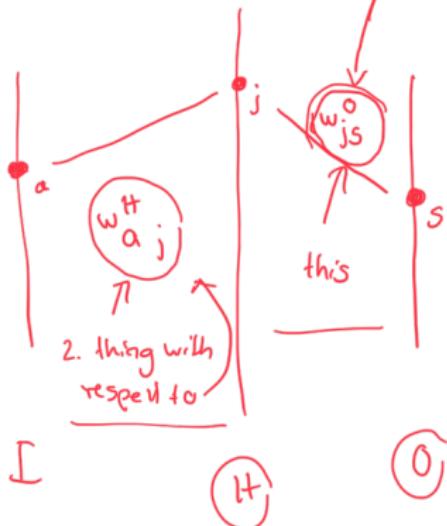
$$w_{ij}^H$$

$$\text{sig}(x_j^H)(1 - \text{sig}(x_j^H) \cdot 1)$$

o

Backpropagation

$$\frac{\partial E}{\partial b_j^4} = \sum_{s=1}^{N_o} (O_s^o - t_s) \cdot \text{sig}(x_s^o) \cdot (1 - \text{sig}(x_s^o)) [w_{j,s}^o] \text{sig}(x_j^h) \cdot (1 - \text{sig}(x_j^h)) \cdot 1$$



Backpropagation

Backpropagation

References i

-  U. M. Ascher and C. Greif.
A First Course on Numerical Methods.
SIAM, Philadelphia, 2011.
-  P. C. Hansen.
Rank-deficient and discrete ill-posed problems.
SIAM Monographs on Mathematical Modeling and
Computation. Society for Industrial and Applied Mathematics
(SIAM), Philadelphia, PA, 1998.

References ii

-  P. C. Hansen.
Discrete inverse problems, volume 7 of Fundamentals of Algorithms.
Society for Industrial and Applied Mathematics (SIAM),
Philadelphia, PA, 2010.
-  J. Nocedal and S. Wright.
Numerical Optimization.
Springer Series in Operations Research and Financial
Engineering. Springer Science & Business Media, New York,
Dec. 2006.

-  Y. Saad.
Iterative Methods for Sparse Linear Systems.
Second Edition. SIAM, Philadelphia, Apr. 2003.

-  C. R. Vogel.
Computational Methods for Inverse Problems.
SIAM, Philadelphia, 2002.