

# 02 - Linear Models

Numerical Methods for Deep Learning

January 10, 2018

# Classification and least-squares regression

Given, examples

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1^\top \\ \mathbf{y}_2^\top \\ \vdots \\ \mathbf{y}_n^\top \end{pmatrix} \in \mathbb{R}^{n \times n_f}$$

and labels

$$\mathbf{C} = \begin{pmatrix} \mathbf{c}_1^\top \\ \mathbf{c}_2^\top \\ \vdots \\ \mathbf{c}_n^\top \end{pmatrix} \in \mathbb{R}^{n \times n_c}$$

Goal: Find a classification/prediction function  $f(\cdot, \boldsymbol{\theta})$ , i.e.,

$$f(\mathbf{y}_j, \boldsymbol{\theta}) \approx \mathbf{c}_j, \quad j = 1, \dots, n.$$

# Regression and least-squares

Simplest option, a linear model with  $\theta = (\mathbf{W}, \mathbf{b})$  and

$$f(\mathbf{Y}, \mathbf{W}, \mathbf{b}) = \mathbf{Y}\mathbf{W} + \mathbf{1}\mathbf{b}^\top \approx \mathbf{C}$$

- ▶  $\mathbf{W} \in \mathbb{R}^{n_f \times n_c}$  are *weights*
- ▶  $\mathbf{b} \in \mathbb{R}^{n_c}$  are *biases*
- ▶  $\mathbf{1} \in \mathbb{R}^n$  is a vector of ones

Equivalent notation:

$$(\mathbf{Y} \quad \mathbf{1}) \begin{pmatrix} \mathbf{W} \\ \mathbf{b}^\top \end{pmatrix} \approx \mathbf{C}$$

Problem may not have a solution, or may have infinite solutions (when?). Solve through optimization

$$\min_{\mathbf{W}} \frac{1}{2} \|\mathbf{Y}\mathbf{W} - \mathbf{C}\|_F^2$$

(Frobenius norm:  $\|\mathbf{A}\|_F^2 = \text{trace}(\mathbf{A}^\top \mathbf{A}) = \sum_{i,j} \mathbf{A}_{i,j}^2$ .)

# Regression and least-squares

To minimize a function need to differentiate and equate to 0

$$\frac{\partial \left( \frac{1}{2} \|\mathbf{Y}\mathbf{W} - \mathbf{C}\|_F^2 \right)}{\partial \mathbf{W}} = 0$$

Computing the derivatives in three steps

1.

$$\frac{\partial \left( \frac{1}{2} \|\mathbf{R}\|_F^2 \right)}{\partial \mathbf{R}} = ???$$

2.

$$\frac{\partial (\mathbf{Y}\mathbf{W})}{\partial \mathbf{W}} = ???$$

3. Use chain rule

# Regression and least-squares

Putting it all together gives

$$\frac{\partial \left( \frac{1}{2} \|\mathbf{Y}\mathbf{W} - \mathbf{C}\|_F^2 \right)}{\partial \mathbf{W}} = \mathbf{Y}^\top (\mathbf{Y}\mathbf{W} - \mathbf{C}) = 0$$

Reorganizing obtain the **normal equations**

$$\mathbf{W} = (\mathbf{Y}^\top \mathbf{Y})^{-1} \mathbf{Y}^\top \mathbf{C}.$$

Assume that  $\mathbf{Y}^\top \mathbf{Y}$  is invertible

- ▶ Sufficient amount of data
- ▶ Data is linearly independent

# Coding: Least-Squares Regression

1. Write a code for solving

$$\min_{\mathbf{W}, \mathbf{b}} \frac{1}{2} \|\mathbf{Y}\mathbf{W} + \mathbf{1}\mathbf{b}^T - \mathbf{C}\|^2$$

and apply it to some of our test data (MNIST / CIFAR10 / Segmentation)

2. Solve the problem using the normal equations derived above.
3. Use optimal weights to predict labels for test data. How good do you get?

## Ill-posedness and regularization - I

If the data is linearly dependent or close to be linearly dependent, least-squares problem gives no good solution. Understanding can be gained by the Singular Value Decomposition (SVD)

$$\mathbf{Y} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$$

where  $\mathbf{U} \in \mathbb{R}^{n \times n_f}$ ,  $\mathbf{V} \in \mathbb{R}^{n_f \times n_f}$  satisfy

$$\mathbf{U}^\top \mathbf{U} = \mathbf{I}, \quad \text{and} \quad \mathbf{V}^\top \mathbf{V} = \mathbf{I}$$

Diagonal of  $\mathbf{\Sigma}$  contains the singular values  $\sigma_1 \geq \dots \sigma_{n_f} \geq 0$

$$\mathbf{\Sigma} = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_{n_f} \end{pmatrix}$$

## Ill-posedness and regularization - I

Important is the effective rank: If  $\sigma_j \ll \sigma_1$  for all  $j \geq k$ , then the effective rank of the problem is  $k$ .

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

Small perturbations in **C** or **Y** yield large perturbations in **W**

Solve regularized problem: For  $\alpha > 0$

$$\min_{\mathbf{W}} \frac{1}{2} \|\mathbf{Y}\mathbf{W} - \mathbf{C}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}\|_F^2$$

*Class assignment - solve the regularized least-squares problem*



## Ill-posedness and regularization - I

Important is the effective rank: If  $\sigma_j \ll \sigma_1$  for all  $j \geq k$ , then the effective rank of the problem is  $k$ .

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

Small perturbations in  $\mathbf{C}$  or  $\mathbf{Y}$  yield large perturbations in  $\mathbf{W}$

Solve regularized problem: For  $\alpha > 0$

$$\min_{\mathbf{W}} \frac{1}{2} \|\mathbf{Y}\mathbf{W} - \mathbf{C}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}\|_F^2$$

*Class assignment - solve the regularized least-squares problem*

$$\mathbf{W} = (\mathbf{Y}^\top \mathbf{Y} + \alpha \mathbf{I})^{-1} \mathbf{Y}^\top \mathbf{C}$$

# The bias variance decomposition

Assume  $\mathbf{C} = \mathbf{Y}\mathbf{W}_{\text{true}} + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \sigma)$ ,  $\alpha > 0$  fixed.

Then setting  $\mathbf{Y}_{\alpha}^{\dagger} = (\mathbf{Y}^{\top}\mathbf{Y} + \alpha\mathbf{I})^{-1}$

$$\begin{aligned}\mathbf{W} - \mathbf{W}_{\text{true}} &= \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\mathbf{C} - \mathbf{W}_{\text{true}} \\ &= (\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\mathbf{Y} - \mathbf{I})\mathbf{W}_{\text{true}} + \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon \\ &= -\alpha\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}} + \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon\end{aligned}$$

# The bias variance decomposition

Assume  $\mathbf{C} = \mathbf{Y}\mathbf{W}_{\text{true}} + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \sigma)$ ,  $\alpha > 0$  fixed.

Then setting  $\mathbf{Y}_{\alpha}^{\dagger} = (\mathbf{Y}^{\top}\mathbf{Y} + \alpha\mathbf{I})^{-1}$

$$\begin{aligned}\mathbf{W} - \mathbf{W}_{\text{true}} &= \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\mathbf{C} - \mathbf{W}_{\text{true}} \\ &= (\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\mathbf{Y} - \mathbf{I})\mathbf{W}_{\text{true}} + \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon \\ &= -\alpha\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}} + \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon\end{aligned}$$



Error depends on  $\epsilon \rightsquigarrow$  take expectation

$$\begin{aligned}\mathbb{E}\|\mathbf{W} - \mathbf{W}_{\text{true}}\|^2 &= \mathbb{E}\|\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon - \alpha\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}}\|^2 \\ &= \underbrace{\alpha^2\|\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}}\|^2}_{\|\text{bias}\|^2} + \underbrace{\sigma^2\text{trace}\left(\mathbf{Y}\mathbf{Y}_{\alpha}^{\dagger\top}\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\right)}_{\text{variance}}\end{aligned}$$

# The bias variance decomposition

Assume  $\mathbf{C} = \mathbf{Y}\mathbf{W}_{\text{true}} + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \sigma)$ ,  $\alpha > 0$  fixed.

Then setting  $\mathbf{Y}_{\alpha}^{\dagger} = (\mathbf{Y}^{\top}\mathbf{Y} + \alpha\mathbf{I})^{-1}$

$$\begin{aligned}\mathbf{W} - \mathbf{W}_{\text{true}} &= \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\mathbf{C} - \mathbf{W}_{\text{true}} \\ &= (\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\mathbf{Y} - \mathbf{I})\mathbf{W}_{\text{true}} + \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon \\ &= -\alpha\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}} + \mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon\end{aligned}$$

Error depends on  $\epsilon \rightsquigarrow$  take expectation

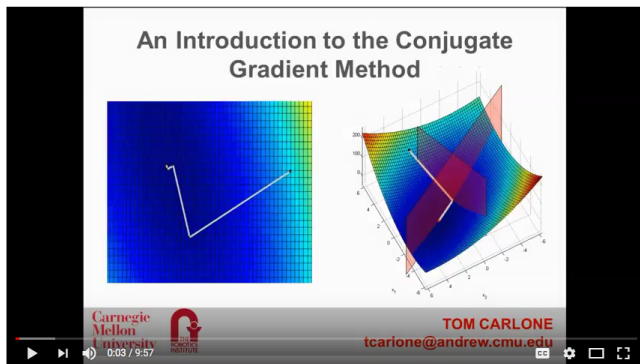
$$\begin{aligned}\mathbb{E}\|\mathbf{W} - \mathbf{W}_{\text{true}}\|^2 &= \mathbb{E}\|\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\epsilon - \alpha\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}}\|^2 \\ &= \underbrace{\alpha^2\|\mathbf{Y}_{\alpha}^{\dagger}\mathbf{W}_{\text{true}}\|^2}_{\|\text{bias}\|^2} + \underbrace{\sigma^2\text{trace}\left(\mathbf{Y}\mathbf{Y}_{\alpha}^{\dagger\top}\mathbf{Y}_{\alpha}^{\dagger}\mathbf{Y}^{\top}\right)}_{\text{variance}}\end{aligned}$$

Take home: No such thing as exact recovery!

# Next time

Solving large-scale least squares problems.

Watch: <https://www.youtube.com/watch?v=eAYohMUpPMA>



Overview of Conjugate Gradient Method