

Linear and Kernel Regression

Linear Models. Regularization. Connection to integral equations.

Machine Learning, Skoltech

Alexey Artemov¹

¹ Skoltech

Practical overview

- › Linear models for regression (a recap)
- › This practical: connection with integral equations
- › Regularization: connection with PCA
- › Trying it out in codes

Multivariate linear regression

- › Multiple features (regressors) $\mathbf{x}_i = (x_{1i}, \dots, x_{di})$ available for each y_i
- › The model:

$$y_1 = w_1 x_{11} + \dots w_d x_{d1} + \varepsilon_1,$$

$$y_2 = w_1 x_{12} + \dots w_d x_{d2} + \varepsilon_2,$$

...

$$y_\ell = w_1 x_{1\ell} + \dots w_d x_{d\ell} + \varepsilon_\ell,$$

is often written in matrix-vector form as

$$\begin{bmatrix} y_1 \\ \vdots \\ y_\ell \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{d1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1\ell} & x_{2\ell} & \dots & x_{d\ell} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_\ell \end{bmatrix} \quad \longleftrightarrow \quad \mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}$$

Multivariate linear regression: the solution

- › The problem: minimize MSE

$$Q(h, X^l) = \sum_{i=1}^{\ell} \left(y_i - \sum_{k=1}^d w_k x_{ki} \right)^2 \equiv \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w} \in \mathbb{R}^d}$$

- › Solve analytically via computing the gradient

$$\nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = 2(\mathbf{y} - \mathbf{X}\mathbf{w})\mathbf{X} = 0$$

- › The solution

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Connection with integral equations

- › The setting (stemming from physics):

$$\mathbf{A}v(y) = u(x),$$

where commonly \mathbf{A} is an integral operator:

$$\mathbf{A}v(y) = \int K(x, y)v(y)dy = u(x),$$

with scalar or vector-valued x, y .

- › Commonly in practice measurements for $u(x)$ taken at points x_1, x_2, \dots, x_ℓ :
 $u_i = u(x_i)$.
- › The problem: estimate $v(y)$ from $u_i, i = 1, \dots, \ell$. Rewriting the integral equation,

$$\sum_{j=1}^n K(x_i, y_j)v_j\Delta y_j = u_i, \quad i = 1, \dots, \ell.$$

Connection with integral equations

- › Letting $\mathbf{u} = (u_1, \dots, u_\ell)$, $\mathbf{v} = (v_1, \dots, v_n)$, $\mathbf{K} = (K(x_i, v_j)\Delta y_j)_{ij}$, we get

$$\mathbf{K}\mathbf{v} = \mathbf{u}.$$

- › In general $n \neq \ell$, this direct inversion of \mathbf{K} is not possible.
- › Commonly to create a square matrix from non-square \mathbf{K} , multiply the equation by \mathbf{K}^\top :

$$\mathbf{K}^\top \mathbf{K} \mathbf{v} = \mathbf{K}^\top \mathbf{u}.$$

Still making $\mathbf{K}^\top \mathbf{K}$ square does not guarantee invertibility (or even a good condition number)!

Ridge regression

- › Ridge regression: adding a “ridge” to the diagonal of $\mathbf{K}^\top \mathbf{K}$ would improve conditioning:

$$(\mathbf{K}^\top \mathbf{K} + \alpha \mathbf{I}) \mathbf{v} = \mathbf{K}^\top \mathbf{u}.$$

- › You all know that solution:

$$\mathbf{v} = (\mathbf{K}^\top \mathbf{K} + \alpha \mathbf{I})^{-1} \mathbf{K}^\top \mathbf{u}.$$

- › Analogously, Kernel regression can be modified from

$$u(x) = \mathbf{r}^\top \mathbf{R}_0^{-1} \mathbf{u} \quad (\mathbf{R}_0)_{ij} = R(x_i, x_j), \quad (\mathbf{r})_i = R(x, x_i)$$

to

$$u(x) = \mathbf{r}^\top (\mathbf{R}_0 + \alpha \mathbf{I})^{-1} \mathbf{u}.$$

Regularizers?



Figure: Large parameter space



Figure: Regularized models

Ad-hoc regularization: motivation

- › Consider the multivariate linear regression problem with $\mathbf{X} \in \mathbb{R}^{d \times d}$

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w} \in \mathbb{R}^d}$$

Ad-hoc regularization: motivation

- › Consider the multivariate linear regression problem with $\mathbf{X} \in \mathbb{R}^{d \times d}$

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w} \in \mathbb{R}^d}$$

- › Analytic solution involves computing the product $\mathbf{R} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

Ad-hoc regularization: motivation

- › Consider the multivariate linear regression problem with $\mathbf{X} \in \mathbb{R}^{d \times d}$

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w} \in \mathbb{R}^d}$$

- › Analytic solution involves computing the product $\mathbf{R} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$
- › If $\mathbf{X} = \text{diag}(\lambda_1, \dots, \lambda_d)$ with $\lambda_1 > \lambda_2 > \dots > \lambda_d \rightarrow 0$
(meaning we're in eigenbasis of \mathbf{X}) then

$$\begin{aligned}\mathbf{R} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \\ &= (\text{diag}(\lambda_1, \dots, \lambda_d) \text{diag}(\lambda_1, \dots, \lambda_d))^{-1} \text{diag}(\lambda_1, \dots, \lambda_d) = \\ &= \text{diag}\left(\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_d}\right), \quad \text{leading to huge diagonal values in } \mathbf{R}\end{aligned}$$

Eigenbasis of X ? Connection with PCA

- › Let $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ be the SVD of \mathbf{X} where $\mathbf{V}^\top\mathbf{V} = \mathbf{I}_d$, $\mathbf{U}\mathbf{U}^\top = \mathbf{U}^\top\mathbf{U} = \mathbf{I}_d$, and \mathbf{S} is a diagonal $d \times d$ matrix. Then:

$$\begin{aligned}\mathbf{R} &= (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top = \\ &= ((\mathbf{U}\mathbf{S}\mathbf{V}^\top)^\top \cdot \mathbf{U}\mathbf{S}\mathbf{V}^\top)^{-1}(\mathbf{U}\mathbf{S}\mathbf{V}^\top)^\top = \\ &= (\mathbf{V}\mathbf{S}\mathbf{U}^\top \cdot \mathbf{U}\mathbf{S}\mathbf{V}^\top)^{-1}\mathbf{V}\mathbf{S}\mathbf{U}^\top = \\ &= (\mathbf{V}\mathbf{S}^2\mathbf{V}^\top)^{-1} \cdot \mathbf{V}\mathbf{S}\mathbf{U}^\top = \\ &= (\mathbf{V}^\top)^{-1}\mathbf{S}^{-2}\mathbf{V}^\top \cdot \mathbf{V}\mathbf{S}\mathbf{U}^\top = \\ &= \mathbf{V}\mathbf{S}^{-2}\mathbf{S}\mathbf{U}^\top = \\ &= \mathbf{V}\tilde{\mathbf{S}}_{\text{linear}}\mathbf{U}^\top\end{aligned}$$

$$(\mathbf{AB})^\top = \mathbf{B}^\top\mathbf{A}^\top$$

$$\mathbf{U}^\top\mathbf{U} = \mathbf{I}_d$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$$

$$\text{orthogonality } \mathbf{V}^\top = \mathbf{V}^{-1}$$

$$\tilde{\mathbf{S}}_{\text{linear}} = \mathbf{S}^{-1}.$$

Eigenbasis of X ? Connection with PCA

- › Doing the same calculations with regularized solution gives:

$$\begin{aligned} \mathbf{R}_{\text{ridge}} &= (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I}_d)^{-1} \mathbf{X}^\top = && \text{same math here} \\ &= \mathbf{V}(\mathbf{S}^2 + \alpha \mathbf{I}_d)^{-1} \mathbf{S} \mathbf{U}^\top = \\ &= \mathbf{V} \tilde{\mathbf{S}}_{\text{ridge}} \mathbf{U}^\top && \tilde{\mathbf{S}}_{\text{ridge}} = (\mathbf{S}^2 + \alpha \mathbf{I}_d)^{-1} \mathbf{S}. \end{aligned}$$

- › Predictions with this kind of operators are given by:

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{R} \mathbf{y} = \mathbf{U} \mathbf{S} \mathbf{V}^\top \cdot \mathbf{V} \tilde{\mathbf{S}}_{\text{method}} \mathbf{U}^\top \cdot \mathbf{y} = \mathbf{U} \tilde{\mathbf{S}} \mathbf{U}^\top \mathbf{y}$$

- › In terms of directions $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_d)$ in data,

$$\hat{\mathbf{y}} = \sum_{j=1}^d \mathbf{u}_j \tilde{S}_{jj} \mathbf{u}_j^\top \mathbf{y}$$

Eigenbasis of \mathbf{X} ? Connection with PCA

- › Values \tilde{S}_{jj} define behavior of predictions:
 - › low (near-zero) values of \tilde{S}_{jj} suppress any influence of \mathbf{u}_j on the prediction;
 - › higher values serve as amplifiers, etc.
- › For the least squares,

$$\tilde{S}_{jj} \equiv 1$$

- › For the ridge regression,

$$\tilde{S}_{jj} \equiv [(\mathbf{S}^2 + \alpha \mathbf{I}_d)^{-1} \mathbf{S}]_{jj} = \frac{\lambda_j^2}{\lambda_j^2 + \alpha}$$

where λ_j are the singular values of \mathbf{X} .

- › If λ_j is small compared to α , then direction \mathbf{u}_j is suppressed.

Why ridge regression works

- › Analytic solution: compute the regularized operator

$$\mathbf{R} = (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top$$

Why ridge regression works

- › Analytic solution: compute the regularized operator

$$\mathbf{R} = (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top$$

- › If $\mathbf{X} = \text{diag}(\lambda_1, \dots, \lambda_d)$ with $\lambda_1 > \lambda_2 > \dots > \lambda_d \rightarrow 0$ (meaning we're in eigenbasis of \mathbf{X}) then

$$\begin{aligned} \mathbf{R} &= (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top = \\ &= (\text{diag}(\lambda_1, \dots, \lambda_d) \text{diag}(\lambda_1, \dots, \lambda_d) + \text{diag}(\alpha, \dots, \alpha))^{-1} \text{diag}(\lambda_1, \dots, \lambda_d) = \\ &= \text{diag}\left(\frac{\lambda_1}{\lambda_1^2 + \alpha}, \dots, \frac{\lambda_d}{\lambda_d^2 + \alpha}\right), \end{aligned}$$

smoothing diagonal values in \mathbf{R}