

Linear Regression: Preliminaries

In the lecture we considered the setup

$$y = \beta^T \mathbf{x} + \varepsilon, \quad (1)$$

with $\mathbb{E}[\varepsilon] = 0$ and $\text{var}[\varepsilon] = \sigma^2$ for some data independent variance. This does not allow for a fixed offset (i.e. the model assumes $y = 0$ for $\mathbf{x} = 0$). To be more flexible, while still keeping the structure we discussed one can proceed as follows. Consider

$$y = \beta_0 + \beta^T \mathbf{x} + \varepsilon = \tilde{\beta}^T \tilde{\mathbf{x}} + \varepsilon, \quad (2)$$

where we have defined $\tilde{\beta} = (\beta_0, \beta^T)^T$ and $\tilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$. Throughout these exercises (and also most often in practice) we will assume to be in this more general case, dropping the tilde from the notation. This means that you always need to remember to include one row of ones at the beginning of your data matrix \mathbf{X} .

1 Plotting of the Residuals (10 pt)

One diagnostic tool for verification that the assumption of a linear model was reasonably correct is the plotting of the residuals. Given our assumptions we consider them to be spread around zero with a constant variation and no discernible pattern.

- i) Without looking at the actual data, fit a linear regression model to `linear-res.npz`¹ and visualize the residuals. Discuss any patterns you observe.
- ii) Extend the data matrix by including also squared features, fit a second model and visualize the residuals again.
- iii) Finally, plot both the data as well as the fitted models.

Solution: See jupyter-notebook.

2 Uncertainty in the Parameter Estimation (10 pt + 2pt)

In the lecture we discussed that given that $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ and assuming the model to be correct, one can show that $\hat{\beta} \sim \mathcal{N}(\beta^*, (\mathbf{X}\mathbf{X}^T)^{-1}\sigma^2)$. In this exercise you will explore this uncertainty for three data sets, `linear-unc1.npz`, `linear-unc2.npz`, and `linear-unc3.npz`. Perform the following steps for each of them. For simplicity we assume the noise variance to be known with $\sigma^2 = 1$ and that we know the intercept/bias² term to be zero (i.e. $\beta_0 = 0$).

- i) Each data set consists of a large number of features \mathbf{x} and target variables y . In order to simulate the uncertainty in $\hat{\beta}$, sample 100 sets of 10 points each from them.³

¹An `.npz` file contains multiple arrays that you can access as you would in a dictionary. You can load it as usual, e.g. `tmp = np.load('foo.npz')` and then access the keys as `list(tmp.keys())`. On this sheet, we will always have the features as `tmp['X']` and the targets as `tmp['Y']`.

²Not the statistical bias, just a naming collision.

³See the last exercise from the first sheet if you need inspiration on how to implement this.

- ii) Fit a linear regression for each subset and plot the values for β_1, β_2 . Discuss the patterns you observe.
- iii) Fit a linear regression to each complete data set and compare your $\hat{\beta}$ estimate it with the results you obtained in ii).
- iv) (technical +2pt) Plot the density of the distribution you obtained for $\hat{\beta}$ in the last step.

Solution: See jupyter-notebook.

3 Estimating Parameter Relevance (10 pt)

Given the importance of linear regression, there exist many approaches to estimate the importance of the individual features after learning the parameters (e.g. hypothesis tests, close inspections of the posteriors in a Bayesian setup, ...). Here you will be implementing an approach that follows from a direct intuitive motivation. In the permutation test to test whether the i -th feature is relevant, you randomly permute it among the data points. Starting from $\mathbf{X} \in \mathbb{R}^{p \times N}$, you then get $\tilde{\mathbf{X}} \in \mathbb{R}^{p \times N}$ with $\tilde{\mathbf{X}}_{i,:} = \pi(\mathbf{X}_{i,:})$, i.e. the i -th row of the data matrix is permuted (where $\pi(\cdot)$ indicates the permutation operation).

You will apply this to `linear-rel.npz`, a variation of the classical boston housing data set. In it, one attempts to estimate the median value of a property given features like the number of rooms, the local crime rate, age of the building, I added one additional random feature that is completely useless prediction-wise and your goal is to find it.

For each of the p features, create a data matrix $\tilde{\mathbf{X}}$ with one of the rows permuted. Fit a new linear regression and compare it to the sum-of-squared residuals of the original \mathbf{X} . If they remain roughly equal then this suggests that the feature is irrelevant. Which feature is the random additional feature?

Solution: See jupyter-notebook.

4 σ^2 Estimation and Heteroscedastic Noise (technical +12pt)

- i) **Maximum Likelihood.** Focusing on a single point (y_n, \mathbf{x}_n) our linear regression model simplifies to

$$y_n = \beta^T \mathbf{x}_n + \varepsilon_n. \quad (3)$$

If we assume that $\varepsilon_n \sim \mathcal{N}(0, \sigma^2)$, this is equivalent to the assumption that $y_n \sim \mathcal{N}(\beta^T \mathbf{x}_n, \sigma^2)$. The logarithm $p(y_n | \beta, \sigma^2)$ is known as the log-likelihood. Having observed N data points this formulation generalizes to a sum of log-likelihoods and we can learn β by maximizing the logarithm of

$$\hat{\beta} = \arg \max_{\beta} \sum_{n=1}^N \log \mathcal{N}(y_n | \beta^T \mathbf{x}_n, \sigma^2). \quad (4)$$

Show that we are solving the same objective as in the SSQ formulation (just with a different scaling factor) and get the same solution for β . This formulation is known as the *Maximum Likelihood* approach, as we learn the parameters that maximize the likelihood of the data.

Solution: We have that

$$\begin{aligned} \sum_{n=1}^N \log \mathcal{N}(y_n | \beta^T \mathbf{x}_n, \sigma^2) &= \sum_{n=1}^N -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma^2) - \frac{(y_n - \beta^T \mathbf{x}_n)^2}{2\sigma^2} \\ &= \underbrace{-\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2)}_{\text{ind of } \beta} - \underbrace{\frac{1}{2\sigma^2}}_{\text{irr factor}} \underbrace{\sum_{n=1}^N (y_n - \beta^T \mathbf{x}_n)^2}_{\text{SSQ}}, \end{aligned}$$

which is again maximized by the $\hat{\beta}$ that minimizes the SSQ.

- ii) **Estimation of σ^2 .** Estimating σ^2 then analogously consists of finding the $\hat{\sigma}^2$ that maximizes this log-likelihood given the estimates $\hat{\beta}$, i.e.

$$\hat{\sigma}^2 = \arg \max_{\sigma^2} \sum_{n=1}^N \log \mathcal{N}(y_n | \hat{\beta}^T \mathbf{x}_n, \sigma^2). \quad (5)$$

Solve this and relate the result to the SSQ residual formulation from the lecture.

Hint: It might be easier to work with the precision $\alpha = \frac{1}{\sigma^2}$ for most of the derivation.

Solution: Given the result of i) and using $\alpha = 1/\sigma^2$ as stated in the hint we have that

$$\sum_{n=1}^N \log \mathcal{N}(y_n | \beta^T \mathbf{x}_n, \sigma^2) = -\frac{N}{2} \log(2\pi) + \frac{N}{2} \log(\alpha) - \frac{\alpha}{2} \text{SSQ}.$$

Setting the gradient with respect to α equal to zero gives us the result that

$$\hat{\sigma}^2 = \frac{1}{\alpha} = \frac{1}{N} \text{SSQ},$$

i.e. that the estimated $\hat{\sigma}^2$ is the residual variance.

- iii) **Heteroscedastic Noise.** The standard formulation of linear regression is of homoscedastic noise, i.e. the variances of the observation noise is independent of \mathbf{x} . A generalization is to have a data point dependent variance on the observation noise, i.e. we have

$$y_n = \beta^T \mathbf{x}_n + \varepsilon_n, \quad (6)$$

with $\mathbb{E}[\varepsilon_n] = 0$ and $\text{var}[\varepsilon_n] = \sigma_n$, which is known as *heteroscedastic noise*. Give the sum-of-squares problem in that case and derive mean and covariance structure of the $\hat{\beta}$ in that case.

Solution: Following our maximum likelihood formulation we see that we have

$$\hat{\beta} = \arg \max \sum_{n=1}^N \frac{1}{\sigma_n^2} (y_n - \beta^T \mathbf{x}_n)^2,$$

i.e. a weighted sum of squares. If we define $\mathbf{S} = \text{diag}(1/\sigma_1^2, \dots, 1/\sigma_N^2)$ the problem becomes in the matrix notation of the lecture

$$\text{SSQ} = (\mathbf{Y} - \beta^T \mathbf{X}) \mathbf{S} (\mathbf{Y} - \beta^T \mathbf{X})^T,$$

and taking the derivatives we analogously get as the solution for β that

$$\hat{\beta} = (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{S} \mathbf{Y}^T.$$

The covariance in turn becomes

$$\begin{aligned} \text{Cov}(\hat{\beta}) &= \text{Cov} \left\{ (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{S} (\mathbf{X}^T \beta + \varepsilon^T) \right\} \\ &= \text{Cov} \left\{ (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{S} \varepsilon^T \right\} \\ &= (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} \mathbf{X} \underbrace{\text{Cov} \{ \varepsilon^T \}}_{=\mathbf{S}^{-1}} \mathbf{X}^T (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} = (\mathbf{X} \mathbf{S} \mathbf{X}^T)^{-1} \end{aligned}$$