

How to choose a model?

Assume that there are k models $M_1 \dots M_k$ that could potentially explain data $(x_1, y_1) \dots (x_n, y_n)$.

What model should we choose?

Bayes formula suggest posterior probabilities

$$P_i[M_i | (x_1, y_1) \dots (x_n, y_n)] = \frac{P_i[(x_1, y_1) \dots (x_n, y_n) | M_i] \cdot P_i[M_i]}{P_i[(x_1, y_1) \dots (x_n, y_n)]}$$

One of these models will have the largest probability. If I must choose only one, I should pick the one with highest probability.

Note that the choice may heavily depend on my prior knowledge encoded as $P_i[M_1] \dots P_i[M_k]$

1

Maximum likelihood estimate for coin flipping

Set of models

Let model M_α denote setting where $x_1 \dots x_n$ are independently drawn from Bernoulli distribution with bias α

$$P_i[x_1 \dots x_n | M_\alpha] = \alpha^k (1-\alpha)^{n-k}$$

where k is the number of ones in the series $x_1 \dots x_n$.

Maximisation task

$$P_i[x_1 \dots x_n | M_\alpha] \xrightarrow{\alpha} \max$$

$$F(\alpha) = \alpha^k (1-\alpha)^{n-k} \xrightarrow{\alpha \in [0,1]} \max$$

Solution $\frac{\partial F}{\partial \alpha} = 0 \Leftrightarrow \alpha^{k-1} (1-\alpha)^{n-k+1} (k(1-\alpha) - (n-k)\alpha) = 0$
 $\Leftrightarrow \alpha = \frac{k}{n}$

3

Maximum likelihood principle

If I am impartial - have no preferences to models:

$$P_i[M_i] = \text{const}, \quad i=1 \dots k$$

then I should always choose the model with highest likelihood

$$P_i[(x_1, y_1) \dots (x_n, y_n) | M_i] \xrightarrow{i} \max.$$

This choice principle is known as maximum likelihood principle. We formulated it in terms of finite number of models. However, we generalise it to infinite number of models:

$$P_i[(x_1, y_1) \dots (x_n, y_n) | M] \xrightarrow{M} \max$$

This can be viewed as a portfolio or a process where we always choose among finite set of models but the precision can be arbitrarily small.

2

Maximum likelihood estimate for lin. regression

Set of models

Let us consider the class of models $M_{a,b}$ in the form x_1, \dots, x_n are drawn from the normal distribution $N(0,1)$ and $y_1 \dots y_n$ are computed as

$$y_i = ax_i + b + \varepsilon_i \quad \text{for } \varepsilon_i \sim N(0, \sigma^2)$$

Now note that under independent draws assumption:

$$P[(x_1, y_1) \dots (x_n, y_n) | a, b] = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_i^2}{2}\right) \cdot \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - ax_i - b)^2}{2\sigma^2}\right)$$
$$= \frac{1}{(2\pi)^n \sigma^n} \cdot \exp\left(-\sum_{i=1}^n \frac{x_i^2}{2}\right) \cdot \underbrace{\exp\left(-\sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{2\sigma^2}\right)}_{F(a,b)}$$

constant w.r.t a, b

4

Further remarks and analysis

Since the density $p(x_i, y_i | a, b)$ contains terms that do not change when we vary parameters a and b we can solve the following maximisation task

$$F(a, b) = \exp\left(-\sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{2\sigma^2}\right) \xrightarrow{a, b} \max$$

Since logarithm is monotonously growing:

$$F(a, b) \rightarrow \max \Leftrightarrow \log F(a, b) \rightarrow \max$$

Hence, we can maximise

$$\sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{2\sigma^2} \xrightarrow{a, b} \max$$

$$\sum_{i=1}^n (y_i - ax_i - b)^2 \xrightarrow{a, b} \min$$

5

Predictions and sanity check

If our assumption about the model is correct then the actual value y that corresponds to x is distributed as

$$y \sim \mathcal{N}(f_w(x), \sigma)$$

So we should output a probability distribution or a 95% confidence interval. For that we need to estimate σ . If we would know the model parameters then $y_i - f_w(x_i)$ would be samples from $\mathcal{N}(0, \sigma)$.

Consequently

$$\sigma \approx \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - f_w(x_i))^2}$$

The latter allows us to draw 95% confidence intervals. You would expect 95% of training samples be inside of these. Secondly you could compare that residuals are from normal distribution with qq-plot.

7

Maximum likelihood and neural networks

Set of models

A fixed network structure with variable weights w . Let $f_w(x)$ denote the output of the network on input x .

Again the distribution of x_i -s is not important. As before, we define

$$y_i = f_w(x_i) + \varepsilon_i \quad \text{for } \varepsilon_i \sim \mathcal{N}(0, \sigma)$$

As a result

$$\log p[(x_i, y_i) | w] = \text{const} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f_w(x_i))^2$$

and we must minimise

$$\sum_{i=1}^n (y_i - f_w(x_i))^2$$

6

Going beyond white Gaussian noise

Normal distribution has light tails. Samples rarely deviate from mean by three or more standard deviations.

If the data contains many outliers, then there are two principle possibilities

- 1) use error distribution with heavier tails
- 2) use explicit modelling of outliers.

Standard solution for the first option is to assume

$$\varepsilon_i \sim \text{Laplace}(0, \beta)$$

Then $\log p(y_i | x_i, \beta) = \text{const} - |y_i - f_w(x_i)|$ and thus the maximum likelihood principle yields

$$\sum_{i=1}^n |y_i - f_w(x_i)| \xrightarrow{w} \min.$$

8

Maximum a posteriori principle

Sometimes we have extra information that makes some models more likely than the others

$$P_i[M_i] \neq \text{const}$$

Then the model with highest likelihood is suboptimal and we should choose the model with highest posterior probability

$$P_i[M_i | (x_1, y_1), \dots, (x_n, y_n)] \xrightarrow{i} \max$$

This method is known as maximum a posteriori principle (MAP)

9

Alternative interpretation to priors

If we take penalised cost function as log of posterior we can reverse the calculation and reverse prior:

$$p[w | \text{Data}] = P[\text{Data} | w] \cdot \exp(-d \cdot \|w\|_1) \cdot c$$

↓

$$p[w] = c_2 \cdot \exp(-d \|w\|_1)$$

and thus we implicitly prefer solutions with low sum of absolute coefficients.

This linear regression algorithm is known as lasso and it leads to solutions with handful non-zero coefficients. These models are used when you need highly interpretable connection between inputs and outputs.

11

Linear regression and MAP

Recall that the restriction

$$\|w\|_1 = |w_0| + \dots + |w_k| \leq c$$

answer that $f(x) = w_0 + w_1 x_1 + \dots + w_k x_k$ is bounded by c in the region $x_i \in [-1, 1]$. Hence, if I know that the output values

$$|f(x_i)| \leq c \quad \text{for all } x_{ij} \in [-1, 1]$$

then I should assign a prior

$$p(w) = \begin{cases} \text{const}, & \|w\|_1 \leq c, \\ 0, & \|w\|_1 > c. \end{cases}$$

Recall that the latter led to minimisation task

$$-\log p((x_1, y_1), \dots, (x_n, y_n) | w) + d \|w\|_1 \xrightarrow{w} \min$$

10

Ridge regression and MAP

Recall that the restriction

$$\|w\|_2^2 = w_0^2 + \dots + w_k^2 \leq c$$

answer that $f(x) = w_0 + w_1 x_1 + \dots + w_k x_k$ is bounded

$$|f(x)| \leq c \quad \text{for all } \|x\|_2^2 \leq 1.$$

Hence, if I want to bound the function in unit ball I should assign a prior

$$p(w) = \begin{cases} \text{const}, & \|w\|_2^2 \leq c \\ 0, & \|w\|_2^2 > c \end{cases}$$

The latter leads to the minimisation task

$$-\log p((x_1, y_1), \dots, (x_n, y_n) | w) + d \|w\|_2^2 \xrightarrow{w} \min$$

12

Linear classification via modelling

Assume that data points come from two distributions.

To be concrete assume that

$$x_i \sim \mathcal{N}(\mu_0, \Sigma_0) \quad \text{if } y_i = 0$$

$$x_i \sim \mathcal{N}(\mu_1, \Sigma_1) \quad \text{if } y_i = 1$$

As a result

$$p(x_1, y_1, \dots, x_n, y_n | \mu_0, \mu_1, \Sigma_0, \Sigma_1) = \prod_{i=1}^n \frac{1}{(2\pi)^{n/2} \det \Sigma_{y_i}} \exp\left(-\frac{1}{2}(x_i - \mu_{y_i})^T \Sigma_{y_i}^{-1} (x_i - \mu_{y_i})\right)$$

Again it makes sense to consider log-likelihood

$$\log p(x_1, y_1, \dots | \mu_0, \mu_1, \Sigma_0, \Sigma_1) = \text{const} + \sum_{y_i=0} -\log \det \Sigma_0 - \frac{1}{2}(x_i - \mu_0)^T \Sigma_0^{-1} (x_i - \mu_0) + \dots$$

13

Unsupervised clustering. The model

When we know the class labels y_i , then it is easy to estimate $\mu_0, \mu_1, \Sigma_0, \Sigma_1$ using maximum likelihood principle and later use Bayes formula to assign probabilities

$$P_0[y=0|x] \quad \text{and} \quad P_1[y=1|x]$$

However, sometimes class labels are missing.

Then we have to update the model

$$P_0[y_i=0] = d_0 \quad P_1[y_i=1] = d_1$$

and estimate the posterior distribution

$$P_0[\mu_0, \mu_1, \Sigma_0, \Sigma_1 | \text{data}]$$

15

Further analysis

It is easy to see that the log-likelihood decomposes into two functions

$$\log p(x_1, y_1, \dots | \Sigma_0, \Sigma_1) = F_0(\Sigma_0, \mu_0) + F_1(\Sigma_1, \mu_1) + \text{const}$$

where

$$F_0(\Sigma_0, \mu_0) = -n_0 \log \Sigma_0 - \frac{1}{2} \sum_{y_i=0} (x_i - \mu_0)^T \Sigma_0^{-1} (x_i - \mu_0)$$

$$F_1(\Sigma_1, \mu_1) = -n_1 \log \Sigma_1 - \frac{1}{2} \sum_{y_i=1} (x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1)$$

where n_0 and n_1 are the counts of labels 0 and 1.

To maximise the log-likelihood we can separately maximise F_0 and F_1 . The latter is equivalent of fitting normal distributions to multivariate data.

14

Analysis in one-dimensional case

Note that

$$p(x_i | \text{model}) = d_0 p(x_i | \mu_0, \sigma_0) + d_1 p(x_i | \mu_1, \sigma_1) \quad i=1 \dots n$$

and thus

$$\log p(x_1, \dots, x_n | \text{model}) = \sum_{i=1}^n \log (d_0 p(x_i | \mu_0, \sigma_0) + d_1 p(x_i | \mu_1, \sigma_1))$$

In principle, we could solve this by computing derivative

$$\frac{\partial}{\partial \mu_0} p, \frac{\partial}{\partial \mu_1} p, \frac{\partial}{\partial \sigma_0} p, \frac{\partial}{\partial \sigma_1} p$$

and then applying some sort of gradient descent algorithm. However, this is analytically difficult.

We see more advanced gradient search algorithms in the following lecture

16