

Machine Learning

Bayesian methods

Wojciech Czech

February 1, 2025

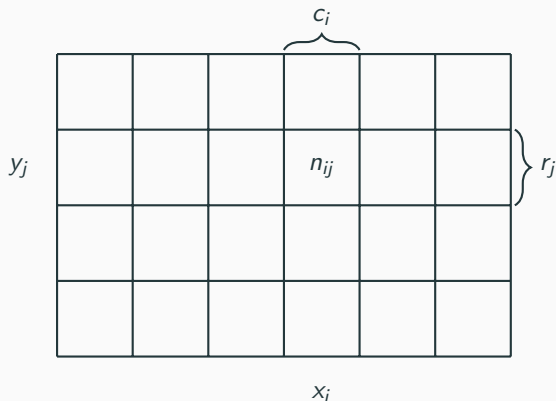
Institute of Computer Science

Table of contents

1. The Rules of Probability
2. Bayes' theorem
3. Bayesian curve fitting
4. Bayesian inference
5. Bayesian linear regression

The Rules of Probability

Sample space of two random variables



Random variable X can take any of the values x_i , where $i = 1, \dots, M$

Random variable Y can take any of the values y_j , where $j = 1, \dots, L$

- N - number of trials, in which we sample both X and Y
- n_{ij} - number of trials, in which we observed $X = x_i$ and $Y = y_j$
- c_i - number of trials, in which we observed $X = x_i$ regardless of value taken by Y
- r_j - number of trials, in which we observed $Y = y_j$ regardless of value taken by X
- We consider limit $N \rightarrow \infty$

Joint probability: $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Joint probability: $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Marginal probability: $p(X = x_i) = \frac{c_i}{N}$

Sum rule

Joint probability: $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Marginal probability: $p(X = x_i) = \frac{c_i}{N}$

$$c_i = \sum_j n_{ij}$$

Sum rule

Joint probability: $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Marginal probability: $p(X = x_i) = \frac{c_i}{N}$

$$c_i = \sum_j n_{ij}$$

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_j n_{ij} = \sum_j \frac{n_{ij}}{N} = \sum_j p(X = x_i, Y = y_j)$$

Sum rule

Joint probability: $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Marginal probability: $p(X = x_i) = \frac{c_i}{N}$

$$c_i = \sum_j n_{ij}$$

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_j n_{ij} = \sum_j \frac{n_{ij}}{N} = \sum_j p(X = x_i, Y = y_j)$$

Sum rule of probability: $p(X = x_i) = \sum_j p(X = x_i, Y = y_j)$

Conditional probability: $p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$

Conditional probability: $p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$

Only trials for which $X = x_i$

Conditional probability: $p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$

Only trials for which $X = x_i$

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \frac{c_i}{N} = p(Y = y_j | X = x_i) p(X = x_i)$$

Conditional probability: $p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$

Only trials for which $X = x_i$

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \frac{c_i}{N} = p(Y = y_j | X = x_i) p(X = x_i)$$

Product rule of probability:

$$p(X = x_i, Y = y_j) = p(Y = y_j | X = x_i) p(X = x_i)$$

Rules of probability

Sum rule

$$p(X) = \sum_Y p(X, Y)$$

Product rule

$$p(X, Y) = p(Y|X)p(X)$$

Symmetry property

$$p(X, Y) = p(Y, X)$$

Example 1

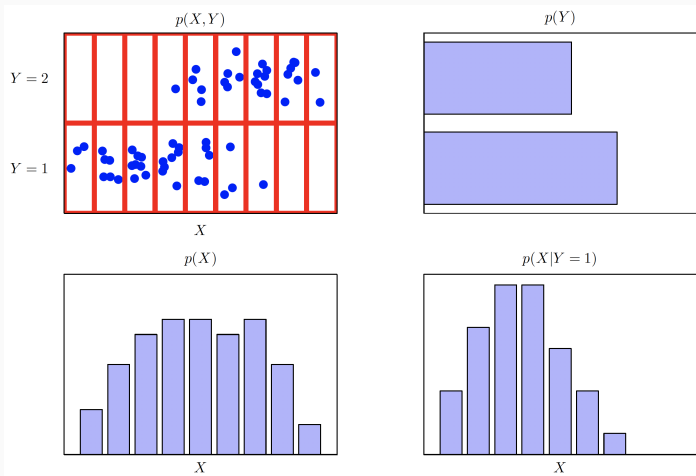


Figure 1: Distribution over two random variables [1].

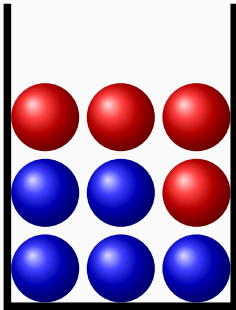
Bayes' theorem

Bayes' theorem

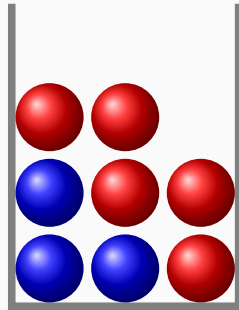
$$p(Y|X) = \frac{p(X, Y)}{p(X)} = \frac{p(Y, X)}{p(X)} = \frac{p(X|Y)p(Y)}{p(X)} \quad (1)$$

$$p(Y|X) = \frac{p(X|Y)p(Y)}{\sum_Y p(Y, X)} = \frac{p(X|Y)p(Y)}{\sum_Y p(X|Y)p(Y)} \quad (2)$$

Example 2: Colored balls



$$p(B = \text{black}) = 6/10$$



$$p(B = \text{gray}) = 4/10$$

Example 2: Probabilities

$$p(C = \text{blue} | B = \text{black}) = 5/9$$

$$p(C = \text{red} | B = \text{black}) = 4/9$$

$$p(C = \text{blue} | B = \text{gray}) = 3/8$$

$$p(C = \text{red} | B = \text{gray}) = 5/8$$

$$\begin{aligned} p(C = \text{red}) &= p(C = \text{red}, B = \text{black}) + p(C = \text{red}, B = \text{gray}) = \\ &= p(C = \text{red} | B = \text{black})p(B = \text{black}) + p(C = \text{red} | B = \text{gray})p(B = \text{gray}) \end{aligned}$$

$$p(C = \text{red}) = 31/60$$

$$p(C = \text{blue}) = 1 - p(C = \text{red}) = 29/60$$

Example 2: Probabilities

$$p(B = \text{black} | C = \text{red}) = \frac{p(C = \text{red} | B = \text{black})p(B = \text{black})}{p(C = \text{red})} = 48/93$$

prior: 0.6, posterior: 0.52

$$p(B = \text{gray} | C = \text{red}) = 1 - p(B = \text{black} | C = \text{red}) = 45/93$$

$$p(B = \text{black} | C = \text{blue}) = \frac{p(C = \text{blue} | B = \text{black})p(B = \text{black})}{p(C = \text{blue})} = 60/87$$

prior: 0.6, posterior: 0.69

$$p(B = \text{gray} | C = \text{blue}) = 1 - p(B = \text{black} | C = \text{blue}) = 27/87$$

Prior and posterior probability

Prior probability

Probability available before observation, e.g., $p(B)$

Posterior probability

Probability obtained after observation, e.g., $p(B|C)$

Bayes' theorem converts a prior probability into a posterior probability by incorporating the evidence provided by the observed data.

Bayesian probability

Classical probability (Frequentist probability)

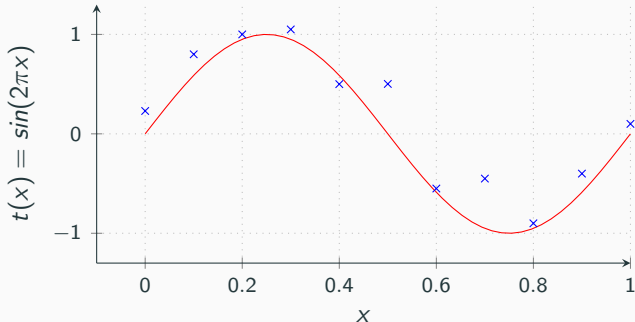
Limit of an event relative frequency in many trials

Bayesian probability

Quantifies uncertainty in the light of new evidence

Bayesian curve fitting

Example 3: Curve fitting

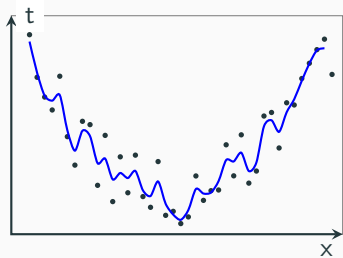
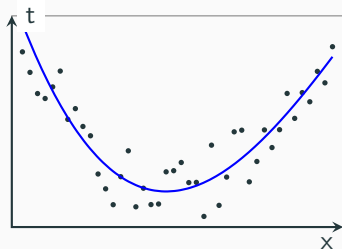
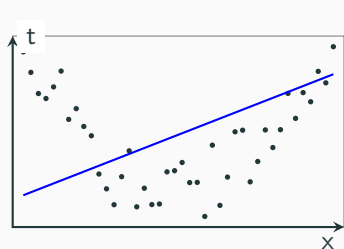


$$\mathbf{x} = [x_1, \dots, x_n]^T, \mathbf{t} = [t_1, \dots, t_n]^T \quad (3)$$

$$y(x, \mathbf{w}) = \sum_{i=0}^m w_i x^i \quad (4)$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=0}^m (y(x_i, \mathbf{w}) - t_i)^2 \quad (5)$$

Model complexity vs. size of training set



Example 3: Curve fitting - regularization

Model complexity and overfitting

- Number of model parameters vs. size of training set
- Values of model parameters

Regularization

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=0}^m (y(x_i, \mathbf{w}) - t_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (6)$$

Example 3: Curve fitting - Bayesian formulation

- Assumptions regarding \mathbf{w} reflected by $p(\mathbf{w})$ - prior probability distribution
- $\mathcal{D} = \{t_1, \dots, t_n\}$ - observed data
- Measure uncertainty in \mathbf{w} after observing \mathcal{D} using posterior probability $p(\mathbf{w}|\mathcal{D})$
- $p(\mathcal{D}|\mathbf{w})$ - likelihood function¹ (function of parameter vector \mathbf{w})

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \quad (7)$$

- Bayes' theorem

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

¹How probable the observed data set is for different settings of \mathbf{w}

Likelihood function

Likelihood function

Measures fit of a statistical model to observed data for a given values of model parameters. Interpreted as a function of parameters only, treating random variables as fixed. Formed from joint probability distribution of the sample.

Maximum likelihood estimation

Finding combination of model parameter values that maximize the probability of drawing the sample obtained.

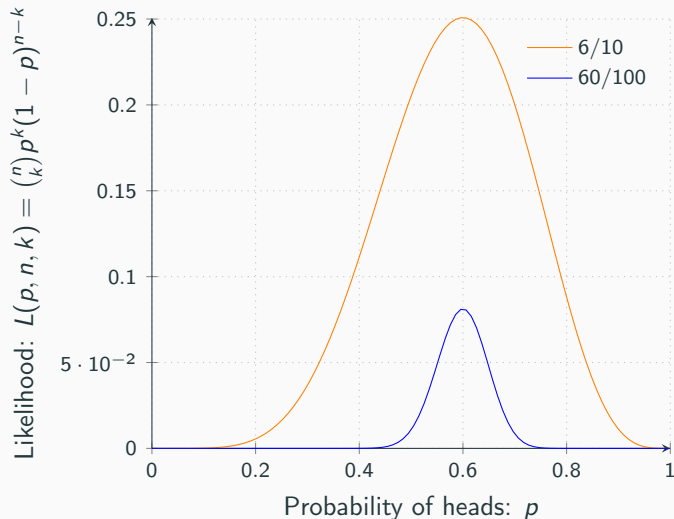
Log likelihood function

Maximizing the likelihood is equivalent to maximizing the log-likelihood and minimizing negative log-likelihood function.

Example 4: Coin flip

- Flipping coin 10 times, obtained 6 heads and 4 tails
- Binomial distribution: $P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$
- H_1 : Fair coin $p(\text{heads}) = w_1 = 0.5$
- $P(X = 6|w_1) \approx 0.21$
- H_2 : Trick coin $p(\text{heads}) = w_2 = 0.75$
- $P(X = 6|w_2) \approx 0.15$
- $L(H_1) = 0.21$, $L(H_2) = 0.15$
- Likelihood ratio: $LR = \frac{L(H_1)}{L(H_2)} = 1.4$
- Hypothesis H_1 explains the data better

Example 4: Coin flip - maximum likelihood



Example 5: Students

Piotr i Paweł spierają się, czy studenci z ich roku preferują Bayesowską (H_1) czy częstotliwościową (H_2) interpretację prawdopodobieństwa. Aby rozstrzygnąć spór, poprosili koleżankę Kasię żeby zapytała 10 losowych studentów opuszczających salę po zajęciach o ich zdanie w spornej kwestii. Dwaj oponenti zgodzili się również wstępnie, że można uznać, że odsetek studentów preferujących interpretację Bayesowską θ podlega rozkładowi $Beta(\theta; \alpha, \beta) = c \theta^{\alpha-1} (1 - \theta)^{\beta-1}$, przy czym $\alpha = 2$ i $\beta = 2$, a c jest stałą normalizującą, którą można obliczyć za pomocą wzoru: $c = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}$, przy czym $\Gamma(m) = (m-1)!$ jeśli m jest liczbą naturalną.

- a) Niech k oznacza liczbę studentów w próbkę zbadanej przez Kasię, preferujących Bayesowską interpretację prawdopodobieństwa. Jaka jest funkcja masy prawdopodobieństwa dla **dyskretnej** zmiennej k ?
- b) Oblicz rozkład *a posteriori* dla θ jeśli $k = 6$.
- c) Znajdź rozwiązanie MAP (maximum a posteriori) dla parametru θ . Czyja hipoteza (H_1 vs. H_2) jest bardziej potwierdzona przez obserwowane dane?

Example 5: solution (1)

- $H_1: \theta > 0.5$ (Piotr)
 - $H_2: \theta \leq 0.5$ (Paweł)
- a) $p(k) = \binom{10}{k} \theta^k (1 - \theta)^{10-k}$ (Kasia odpytuje wychodzących)
- b) $p(\theta) = c \theta(1 - \theta)$ (**prior**, zgodzili się wstępnie)
- $p(k = 6|\theta) = \binom{10}{6} \theta^6 (1 - \theta)^4$ (**likelihood**)
- $p(\theta|k = 6) \propto \binom{10}{6} \theta^6 (1 - \theta)^4 c \theta(1 - \theta) \propto \theta^7 (1 - \theta)^5$ (**posterior**,
nieznormalizowana)

Można zauważyć, że rozkład posterior jest rozkładem $Beta(\theta; 8, 6)$ czyli, że: $\alpha = 8$, $\beta = 6$.

Zamiast całkować, skorzystamy ze wzoru: $c = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} = \frac{13!}{7!5!} = 10296$

$$p(\theta|k = 6) = \frac{13!}{7!5!} \theta^7 (1 - \theta)^5$$

Example 5: solution (2)

c)

$$\frac{\partial}{\partial \theta} \theta^7 (1 - \theta)^5 = 0$$

$$\theta^6 (1 - \theta)^4 (7 - 12\theta) = 0$$

$$\theta = 0 \text{ lub } \theta = \frac{7}{12}$$

$$\text{Maximum dla } \theta = \frac{7}{12} \approx 0.58$$

Dane bardziej wspierają hipotezę Piotra

Bayesian curve fitting

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \quad (8)$$

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \quad (9)$$

Precision: $\beta^{-1} = \sigma^2$

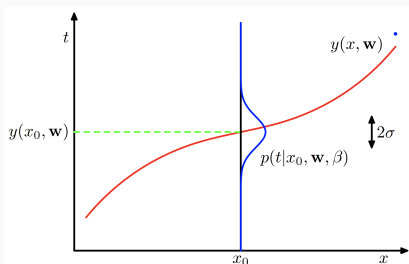


Figure 2: See: [1].

Bayesian curve fitting - model parameters

$$\mathbf{x} = [x_1, \dots, x_n]^T, \mathbf{t} = [t_1, \dots, t_n]^T \quad (10)$$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{i=0}^n \mathcal{N}(t_i|y(x_i, \mathbf{w}), \beta^{-1}) \quad (11)$$

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{i=0}^n (y(x_i, \mathbf{w}) - t_i)^2 + \frac{n}{2} \ln \beta - \frac{n}{2} \ln (2\pi) \quad (12)$$

- Data drawn independently from distribution Eq. 8 (i.i.d.)
- Maximize likelihood with respect to \mathbf{w} - maximize log likelihood - minimize negative log likelihood - minimize root square error function (see: Equation 5).
- Sum-of-squares error function appears as a consequence of maximizing likelihood under the assumption of a Gaussian noise distribution.

Bayesian curve fitting - maximum likelihood

Maximum likelihood parameters

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=0}^n (y(x_i, \mathbf{w}) - t_i)^2 \quad (13)$$

$$\frac{1}{\beta_{ML}} = \sum_{i=1}^n (y(x_i, \mathbf{w}_{ML}) - t_i)^2 \quad (14)$$

Predictive distribution

$$p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t|y(x, \mathbf{w}_{ML}), \beta_{ML}^{-1}) \quad (15)$$

Bayesian curve fitting - maximum posterior

Prior distribution over polynomial coefficients

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(m+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\} \quad (16)$$

Bayes' theorem

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha) \quad (17)$$

MAP - maximum posterior

Determine most probable value of \mathbf{w} given data

$$\mathbf{w}_{MP} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{\beta}{2} \sum_{i=0}^n (y(x_i, \mathbf{w}) - t_i)^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\} \quad (18)$$

See Equation 6 and Equation 18 : $\lambda = \alpha/\beta$

Posterior predictive distribution

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t})d\mathbf{w} \quad (19)$$

$p(t|x, \mathbf{w})$ - see Equation 8

$p(\mathbf{w}|\mathbf{x}, \mathbf{t})$ - posterior distribution over parameters, see Equation 17

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x)) \quad (20)$$

$$m(x) = \beta\phi(x)^T \mathbf{S} \sum_{i=0}^n \phi(x_i) t_i \quad (21)$$

$$s^2(x) = \beta^{-1} + \phi(x)^T \mathbf{S} \phi(x) \quad (22)$$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{i=0}^n \phi(x_i) \phi(x_i)^T \quad (23)$$

$$\phi(x) = [1, x, \dots, x^m]^T \quad (24)$$

Bayesian inference

Example 6: Binary classification problem

- \mathbf{x} - input feature vector
- t - output target variable: $t \in \{\mathcal{C}_1, \mathcal{C}_2\}$, alternatively $t \in [0, 1]$
- Inference problem - determining $p(\mathbf{x}, t)$
- Inference problem vs. decision problem
- Intuition: select the class having the higher posterior probability

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \quad (25)$$

Minimizing misclassification rate

Decision region \mathcal{R}_k

All points \mathbf{x} belonging to \mathcal{R}_k are assigned to class \mathcal{C}_k

Decision boundaries

The boundaries between decision regions

$$p(\text{mistake}) = p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \quad (26)$$

$$p(\text{mistake}) = \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x} \quad (27)$$

Decision rule

If $p(\mathbf{x}, \mathcal{C}_1) > p(\mathbf{x}, \mathcal{C}_2)$, then assign \mathbf{x} to class \mathcal{C}_1

If $p(\mathcal{C}_1|\mathbf{x}) > p(\mathcal{C}_2|\mathbf{x})$, then assign \mathbf{x} to class \mathcal{C}_1

Select largest posterior probability to minimize probability of making a mistake

Decision boundary

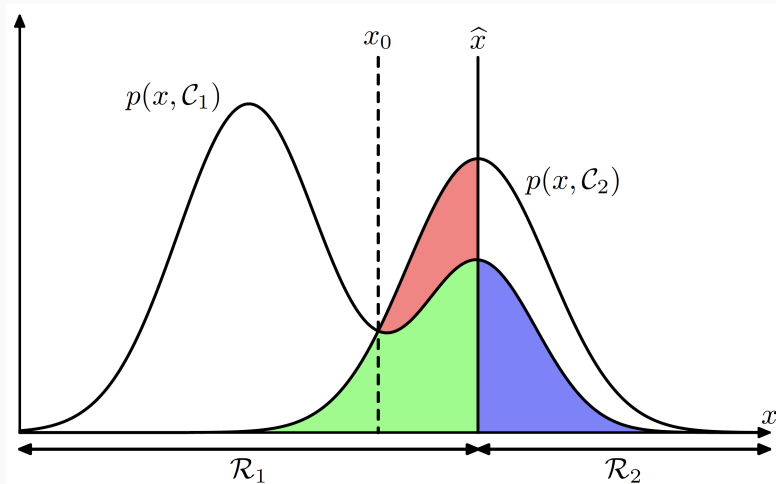


Figure 3: Decision boundary [1].

Minimizing the expected loss

Loss (cost) function

Overall measure of loss incurred in taking any of the available decisions or actions

Loss matrix L

L_{kj} - the cost of assigning \mathbf{x} to class j , while the true class is k

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, C_k) d\mathbf{x} \quad (28)$$

Decision rule

Assign new \mathbf{x} to class j for which the quantity $\sum_k L_{kj} p(C_k|\mathbf{x})$ is a minimum

Classification problem: inference and decision

Inference stage

Use training data to learn a model for $p(\mathcal{C}_k|\mathbf{x})$

Decision stage

Use posterior probabilities to assign classes

Discriminat function

Solve inference and decision together by learning discriminant function, which maps \mathbf{x} directly to decisions

Generative models

- Determine class conditional probability densities $p(\mathbf{x}|\mathcal{C}_k)$ for each class \mathcal{C}_k
- Determine class prior probabilities $p(\mathcal{C}_k)$
- Use Bayes' theorem to find posterior probabilities $p(\mathcal{C}_k|\mathbf{x})$
- Use decision theory to determine class membership for each new input \mathbf{x}
- Sampling model can generate new synthetic data

- Determine posterior probabilities $p(\mathcal{C}_k|\mathbf{x})$ directly
- Use decision theory to determine class membership for each new input \mathbf{x}

- Find function $f(\mathbf{x})$, which maps input vector \mathbf{x} directly to class label
- $f : \mathbf{x} \rightarrow \{\mathcal{C}_1, \dots, \mathcal{C}_k\}$

Advantages of knowing posterior probabilities²

Minimizing risk

Easy revision of loss matrix

Reject option

Exclude regions with big uncertainty about class membership. Simpler adjustment of reject criteria.

Compensating for class priors

Adjusting to imbalanced training datasets

Combining models

Utilizing conditional independence property

²vs. knowing discriminative function only

Naive Bayes Classifier

- Divide classification problem into subproblems with different input spaces
- Obtain posterior probabilities for the classes for each input space separately
- Combine output using rules of probability (normalization)
- **Assume that inputs are independent (conditional independence)**

$$p(\mathbf{x}_I, \mathbf{x}_B | \mathcal{C}_k) = p(\mathbf{x}_I | \mathcal{C}_k) p(\mathbf{x}_B | \mathcal{C}_k) \quad (29)$$

$$\begin{aligned} p(\mathcal{C}_k | \mathbf{x}_I, \mathbf{x}_B) &\propto p(\mathbf{x}_I, \mathbf{x}_B | \mathcal{C}_k) p(\mathcal{C}_k) \propto p(\mathbf{x}_I | \mathcal{C}_k) p(\mathbf{x}_B | \mathcal{C}_k) p(\mathcal{C}_k) \quad (30) \\ &\propto \frac{p(\mathcal{C}_k | \mathbf{x}_I) p(\mathcal{C}_k | \mathbf{x}_B)}{p(\mathcal{C}_k)} \end{aligned}$$

Naive Bayes Classifier - multiple features

$$p(C_k|x_1, \dots, x_n) = \frac{p(x_1|C_k)p(x_2|C_k) \dots p(x_n|C_k)p(C_k)}{p(x_1)p(x_2) \dots p(x_n)} \quad (31)$$

$$p(C_k|x_1, \dots, x_n) = \frac{p(C_k) \prod_{i=1}^n p(x_i|C_k)}{\prod_{i=1}^n p(x_i)} \quad (32)$$

$$p(C_k|x_1, \dots, x_n) \propto p(C_k) \prod_{i=1}^n p(x_i|C_k) \quad (33)$$

$$C = \operatorname{argmax}_{C_k} \left\{ p(C_k) \prod_{i=1}^n p(x_i|C_k) \right\} \quad (34)$$

Example 7: Weather conditions for playing golf

	Outlook x_1	Temperature x_2	Humidity x_3	Windy x_4	Play Golf
1	Rainy	Hot	High	False	No
2	Rainy	Hot	High	True	No
3	Overcast	Hot	High	False	Yes
4	Sunny	Mild	High	False	Yes
5	Sunny	Cool	Normal	False	Yes
6	Sunny	Cool	Normal	True	No
7	Overcast	Cool	Normal	True	Yes
8	Rainy	Mild	High	False	No
9	Rainy	Cool	Normal	False	Yes
10	Sunny	Mild	Normal	False	Yes
11	Rainy	Mild	Normal	True	Yes
12	Overcast	Mild	High	True	Yes
13	Overcast	Hot	Normal	False	Yes
14	Sunny	Mild	High	True	No

Example 7: Prior

	No. cases	$p(\mathcal{C}_k)$
Yes (\mathcal{C}_1)	9	9/14
No (\mathcal{C}_2)	5	5/14
Total	14	1

Example 7: Likelihood

	Yes (C_1)	No (C_2)	$p(x_1 C_1)$	$p(x_1 C_2)$
Sunny	3	2	1/3	2/5
Overcast	4	0	4/9	0
Rainy	2	3	2/9	3/5
Total	9	5	1	1

	Yes (C_1)	No (C_2)	$p(x_2 C_1)$	$p(x_2 C_2)$
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cool	3	1	1/3	1/5
Total	9	5	1	1

	Yes (C_1)	No (C_2)	$p(x_3 C_1)$	$p(x_3 C_2)$
High	3	4	1/3	4/5
Normal	6	1	2/3	1/5
Total	9	5	1	1

	Yes (C_1)	No (C_2)	$p(x_4 C_1)$	$p(x_4 C_2)$
False	6	2	2/3	2/5
True	3	3	1/3	3/5
Total	9	5	1	1

Example 7: Posterior

Query: $\mathbf{x}_q = [x_1, x_2, x_3, x_4] = [\textit{Sunny}, \textit{Hot}, \textit{Normal}, \textit{False}]$

$$p(C_1|\mathbf{x}_q) \propto \frac{1}{3} \frac{2}{9} \frac{2}{3} \frac{2}{3} \frac{9}{14} = 0.02116$$

$$p(C_2|\mathbf{x}_q) \propto \frac{2}{5} \frac{2}{5} \frac{1}{5} \frac{2}{5} \frac{5}{14} = 0.00457$$

$$p(C_1|\mathbf{x}_q) = \frac{0.02116}{0.02116+0.00457} = 0.822$$

$$p(C_2|\mathbf{x}_q) = 1 - 0.822 = 0.178$$

Prediction C_1 : Golf will be played

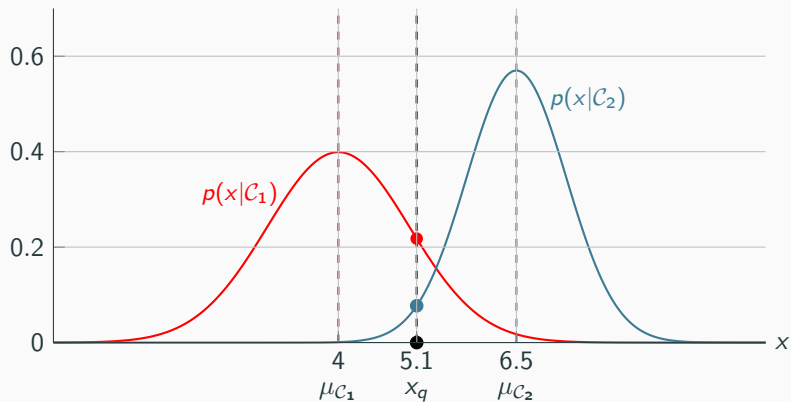
Gaussian Naive Bayes classifier

- x_i has continuous values
- Assume that $p(x_i|\mathcal{C}_k)$ is Gaussian

$$p(x_i|\mathcal{C}_k) = \frac{1}{\sqrt{2\pi\sigma_{\mathcal{C}_k}^2}} \exp \left\{ -\frac{(x_i - \mu_{\mathcal{C}_k})^2}{2\sigma_{\mathcal{C}_k}^2} \right\} \quad (35)$$

- $\mu_{\mathcal{C}_k}, \sigma_{\mathcal{C}_k}^2$ - mean and variance of continuous x_i calculated for a given class \mathcal{C}_k

Gaussian Naive Bayes classifier - visualization



Example 8: Classify gender of the person

	Height x_1	Weight x_2	Foot size x_3	Gender
1	6	180	12	Male
2	5.92	190	11	Male
3	5.58	170	12	Male
4	5.92	165	10	Male
5	5	100	6	Female
6	5.5	150	8	Female
7	5.42	130	7	Female
8	5.75	150	9	Female

Example 8: Means and variances

Gender	μ_{x_1}	$\sigma_{x_1}^2$	μ_{x_2}	$\sigma_{x_2}^2$	μ_{x_3}	$\sigma_{x_3}^2$
Male	5.885	3.5033×10^{-2}	176.25	1.2292×10^2	11.25	9.1667×10^{-1}
Female	5.415	9.7225×10^{-2}	132.5	5.5833×10^2	7.5	1.6667

Male : \mathcal{C}_1 , Female : \mathcal{C}_2

$p(\mathcal{C}_1) = 0.5, p(\mathcal{C}_2) = 0.5$

Query: $\mathbf{x}_q = [x_1, x_2, x_3] = [6, 130, 8]$

Example 8: Likelihood \mathcal{C}_1

$$p(\mathcal{C}_1|\mathbf{x}_q) \propto p(\mathcal{C}_1) \prod_{i=1}^3 p(x_i|\mathcal{C}_1)$$

$$p(x_1 = 6|\mathcal{C}_1) = \frac{1}{\sqrt{2\pi\sigma_{x_1}^2}} \exp\left\{-\frac{(6 - \mu_{x_1})^2}{2\sigma_{x_1}^2}\right\} \approx 1.57889$$

$$p(x_2 = 130|\mathcal{C}_1) = \frac{1}{\sqrt{2\pi\sigma_{x_2}^2}} \exp\left\{-\frac{(130 - \mu_{x_2})^2}{2\sigma_{x_2}^2}\right\} \approx 5.9881 \times 10^{-6}$$

$$p(x_3 = 8|\mathcal{C}_1) = \frac{1}{\sqrt{2\pi\sigma_{x_3}^2}} \exp\left\{-\frac{(8 - \mu_{x_3})^2}{2\sigma_{x_3}^2}\right\} \approx 1.3112 \times 10^{-3}$$

$$p(\mathcal{C}_1|\mathbf{x}_q) \propto 6.1984 \times 10^{-9}$$

Example 8: Likelihood \mathcal{C}_2

$$p(\mathcal{C}_2|\mathbf{x}_q) \propto p(\mathcal{C}_2) \prod_{i=1}^3 p(x_i|\mathcal{C}_2)$$

$$p(x_1 = 6|\mathcal{C}_2) \approx 2.2346 \times 10^{-1}$$

$$p(x_2 = 130|\mathcal{C}_2) \approx 1.6789 \times 10^{-2}$$

$$p(x_3 = 8|\mathcal{C}_2) \approx 2.8669 \times 10^{-1}$$

$$p(\mathcal{C}_2|\mathbf{x}_q) \propto 5.3778 \times 10^{-4}$$

Example 8: Posterior

$$p(\mathcal{C}_1|\mathbf{x}_q) = \frac{6.1984 \times 10^{-9}}{6.1984 \times 10^{-9} + 5.3778 \times 10^{-4}} = 1.1526 \times 10^{-5}$$
$$p(\mathcal{C}_2|\mathbf{x}_q) = 1 - 1.1526 \times 10^{-5} \approx 1$$

Prediction \mathcal{C}_2 : Female

Naive Bayes classifier - advantages

- If the assumption about conditional independence of inputs is true - performs better than other models
- Requires relatively small sizes of training data, fast training
- Easy to implement
- Each distribution can be independently estimated as a one-dimensional distribution (avoiding problems with curse of dimensionality)
- Suitable for multiclass classification
- Simple update of probabilities after new data is received
- Ability to use more complex distributions for continuous input

Naive Bayes classifier - disadvantages

- Very strong assumption of conditional independence of inputs
- Zero frequency problem
- Weights equally all input features
- Problems with handling countinuous variables (discretization vs. assumptions regarding probability distribution)
- Problems with handling imbalanced datasets
- Bad estimator of probability (independence assumption is false if you are trying to estimate the probability)

Naive Bayes classifier - hints

- Reduce highly correlated features (e.g. using PCA)
- Execute domain specific feature engineering
- Try introducing more realistic prior probabilities (domain knowledge)
- Use ensemble methods: bagging or boosting to reduce variance (improve generalization capabilities)
- Use Box-Cox or Yeo-Johnson transformations to make features more Gaussian
- Use **Laplace smoothing** to handle records with zero values in x_i dimension

Naive Bayes classifier - applications

- Real time prediction
- Multi-class predictions
- Text classification
- Spam filtering
- Sentiment analysis
- Recommendation systems

Multivariate Gaussian distribution

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (36)$$

$\mathbf{x} \in \mathbb{R}^{D \times 1}$

$\boldsymbol{\mu} \in \mathbb{R}^{D \times 1}$ - mean vector

$\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$ - covariance matrix

$|\boldsymbol{\Sigma}|$ - determinant of $\boldsymbol{\Sigma}$

$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$ - Mahalanobis distance

Multivariate Gaussian distribution - properties

- Constant on surfaces in \mathbf{x} -space for which Mahalanobis distance is constant
- Elliptical surface of constant probability density for a 2D Gaussian
- Eigendecomposition of covariance matrix

$$\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T \quad (37)$$

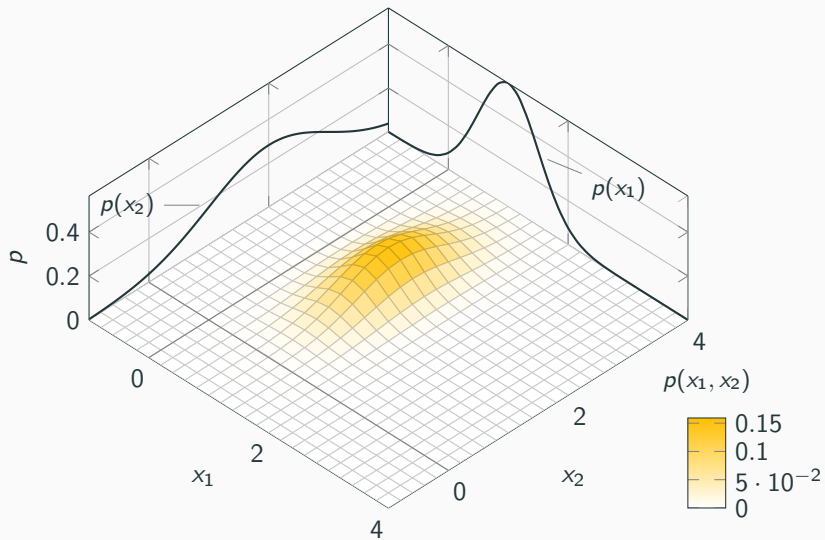
$$\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (38)$$

$$|\Sigma|^{1/2} = \prod_{i=1}^D \lambda_i^{1/2} \quad (39)$$

Multivariate Gaussian distribution - properties

- Number of free parameters $D(D + 3)/2$
- Covariance matrix restricted to diagonal: $2D$ parameters
- Covariance matrix restricted to scaled identity $D + 1$ parameters
- Unimodal distribution (single maximum)

Example 9



Partitioned Gaussian distributions

Two sets of variables

If two sets of variables are jointly Gaussian, then conditional distribution of one set conditioned on the other is also Gaussian. Marginal distribution of either set is also Gaussian.

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix} \quad (40)$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}, \boldsymbol{\Sigma} = \boldsymbol{\Sigma}^T, \boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^T \quad (41)$$

$$\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}, \boldsymbol{\Lambda} = \begin{bmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{bmatrix} \quad (42)$$

Conditional distribution

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad (43)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (44)$$

Marginal distribution

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}) \quad (45)$$

Bayesian linear regression

Linear basis function model

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^{m-1} w_i \phi_i(\mathbf{x}) \quad (46)$$

$\phi_i(\mathbf{x})$ - basis function, w_0 - bias

$$\phi_0(\mathbf{x}) = 1$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{i=0}^{m-1} w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (47)$$

$$\mathbf{w} = [w_0, \dots, w_{m-1}]^T$$

$$\boldsymbol{\phi} = [\phi_0, \dots, \phi_{m-1}]^T$$

Non-linear basis functions

- Polynomial: $\phi_i(x) = x^i$
- Gaussian: $\phi_i(x) = \exp \left\{ -\frac{(x-\mu_i)^2}{2s^2} \right\}$
- Sigmoidal: $\phi_i(x) = \sigma \left(\frac{x-\mu_i}{s} \right)$, $\sigma(a) = \frac{1}{1+\exp(-a)}$
- Fourier
- Wavelet

Bayesian linear basis function model - formulation

$$t(\mathbf{x}) = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad (48)$$

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) = \mathcal{N}(t|\mathbf{w}^T \phi(\mathbf{x}), \beta^{-1}) \quad (49)$$

Bayesian linear basis function model - likelihood

$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ - set of multidimensional inputs

$$\mathbf{t} = [t_1, \dots, t_n]^T \quad (50)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^n \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}), \beta^{-1}) \quad (51)$$

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{i=0}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}))^2 + \frac{n}{2} \ln \beta - \frac{n}{2} \ln (2\pi) \quad (52)$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=0}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}))^2 \quad (53)$$

Maximum likelihood parameters

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = 0 \quad (54)$$

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (55)$$

$$\frac{1}{\beta_{ML}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}_{ML}^T \phi(\mathbf{x}_i) - t_i)^2 \quad (56)$$

$$\Phi = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{m-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_{m-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_n) & \phi_1(\mathbf{x}_n) & \dots & \phi_{m-1}(\mathbf{x}_n) \end{bmatrix} \quad (57)$$

Bayesian linear basis function model - maximum posterior

Prior distribution over model parameters

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \quad (58)$$

Posterior distribution

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_n, \mathbf{S}_n) \quad (59)$$

$$\mathbf{m}_n = \mathbf{S}_n(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t}) \quad (60)$$

$$\mathbf{S}_n^{-1} = \mathbf{S}_0^{-1} + \beta\Phi^T\Phi \quad (61)$$

MAP - maximum posterior

Determine most probable value of \mathbf{w} given data (mode equals mean)

$$\mathbf{w}_{MP} = \mathbf{m}_n \quad (62)$$

Zero-mean isotropic Gaussian prior

Prior distribution

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathbf{I}) \quad (63)$$

Posterior distribution

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_n, \mathbf{S}_n) \quad (64)$$

$$\mathbf{m}_n = \beta \mathbf{S}_n \Phi^T \mathbf{t} \quad (65)$$

$$\mathbf{S}_n^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi \quad (66)$$

MAP - maximum posterior

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{i=1}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const} \quad (67)$$

Online Bayesian regression

- Sequence of :
 - Feature vectors: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots$
 - Target variables $t_1, t_2, \dots, t_i, \dots$
 - Model parameters: $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_i, \dots$

$$\mathbf{w}_i \xrightarrow{t_i} \mathbf{w}_{i+1}$$

- Current prior distribution of model parameters $p(\mathbf{w}_i)$
- **Predictive distribution** $p(t_i|\mathbf{x}_i)$
- Likelihood $p(t_i|\mathbf{x}_i, \mathbf{w}_i)$
- Posterior distribution $p(\mathbf{w}_{i+1}|\mathbf{w}_i, \mathbf{x}_i, t_i)$

$$p(\mathbf{w}_{i+1}|\mathbf{w}_i, \mathbf{x}_i, t_i) \propto p(t_i|\mathbf{x}_i, \mathbf{w}_i)p(\mathbf{w}_i) \quad (68)$$

$$p(t_i|\mathbf{x}_i) = \int p(t_i|\mathbf{x}_i, \mathbf{w})p(\mathbf{w})d\mathbf{w} \quad (69)$$

$$p(t_i|\mathbf{x}_i) \propto p(t_i|\mathbf{x}_i, \mathbf{w}_i)p(\mathbf{w}_i) \quad (70)$$

From Eq. 70 and Eq. 68, assuming initial prior $p(\mathbf{w}_0)$:

$$p(t_0|\mathbf{x}_0) \propto p(t_0|\mathbf{x}_0, \mathbf{w}_0)p(\mathbf{w}_0) \quad (71)$$

After first observation:

$$p(\mathbf{w}_1|\mathbf{w}_0, \mathbf{x}_0, t_0) \propto p(t_0|\mathbf{x}_0, \mathbf{w}_0)p(\mathbf{w}_0) \quad (72)$$

Predictive distribution:

$$p(t_1|\mathbf{x}_1) \propto p(t_1|\mathbf{x}_1, \mathbf{w}_1) \underbrace{p(t_0|\mathbf{x}_0, \mathbf{w}_0)p(\mathbf{w}_0)}_{p(\mathbf{w}_1)} \quad (73)$$

The prior of the weights for the current iteration is the posterior of the weights at the previous iteration.

After second observation:

$$p(\mathbf{w}_2 | \mathbf{w}_1, \mathbf{x}_1, t_1) \propto p(t_1 | \mathbf{x}_1, \mathbf{w}_1) \underbrace{p(t_0 | \mathbf{x}_0, \mathbf{w}_0) p(\mathbf{w}_0)}_{p(\mathbf{w}_1)} \quad (74)$$

After third observation:

$$p(\mathbf{w}_3 | \mathbf{w}_2, \mathbf{x}_2, t_2) \propto p(t_2 | \mathbf{x}_2, \mathbf{w}_2) \underbrace{p(t_1 | \mathbf{x}_1, \mathbf{w}_1) \underbrace{p(t_0 | \mathbf{x}_0, \mathbf{w}_0) p(\mathbf{w}_0)}_{p(\mathbf{w}_1)}}_{p(\mathbf{w}_2)} \quad (75)$$

- Model: $t_i = \mathbf{w}_i^T \mathbf{x}_i + \epsilon_i$
- Likelihood: $p(t_i | \mathbf{x}_i, \mathbf{w}_i) = \mathcal{N}(t_i | \mathbf{w}_i^T \mathbf{x}_i, \beta^{-1})$
- Precision: $\beta = \frac{1}{\sigma^2}$
- Prior: $p(\mathbf{w}_0) = \mathcal{N}(\mathbf{w}_0 | \mathbf{m}_0, \mathbf{S}_0)$
- Mean: $\mathbf{m}_0 = [0, \dots, 0]$
- Covariance: $\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$

$$p(\mathbf{w}_{i+1} | \mathbf{w}_i, \mathbf{x}_i, t_i) = \mathcal{N}(\mathbf{w}_{i+1} | \mathbf{m}_{i+1}, \mathbf{S}_{i+1}) \quad (76)$$

$$\mathbf{S}_{i+1} = (\mathbf{S}_i^{-1} + \beta \mathbf{x}_i \mathbf{x}_i^T)^{-1} \quad (77)$$

$$\mathbf{m}_{i+1} = \mathbf{S}_{i+1}(\mathbf{S}_i^{-1} \mathbf{m}_i + \beta t_i \mathbf{x}_i) \quad (78)$$

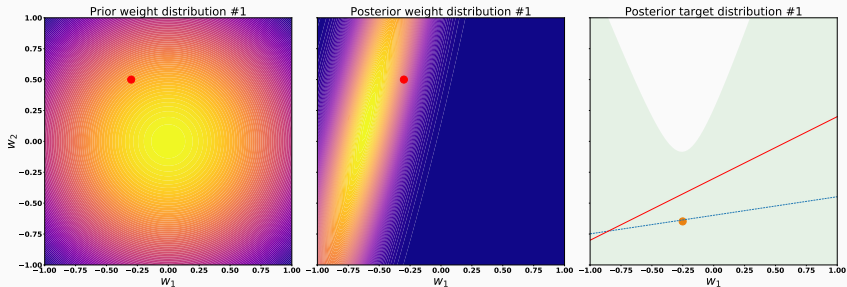
$$p(t_i|\mathbf{x}_i) = \mathcal{N}(t_i|\mu_i, \sigma_i) \quad (79)$$

$$\mu_i = \mathbf{w}_i^T \mathbf{x}_i \quad (80)$$

$$\sigma_i = \frac{1}{\beta} + \mathbf{x}_i^T \mathbf{S}_i \mathbf{x}_i \quad (81)$$

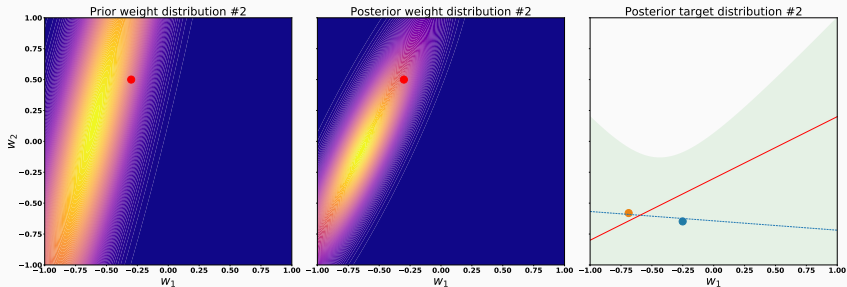
Online linear regression - progressive validation

$$y = w_1 + w_2x, \quad t = -0.3 + 0.5x + \epsilon, \quad \beta = 25, \alpha = 2$$



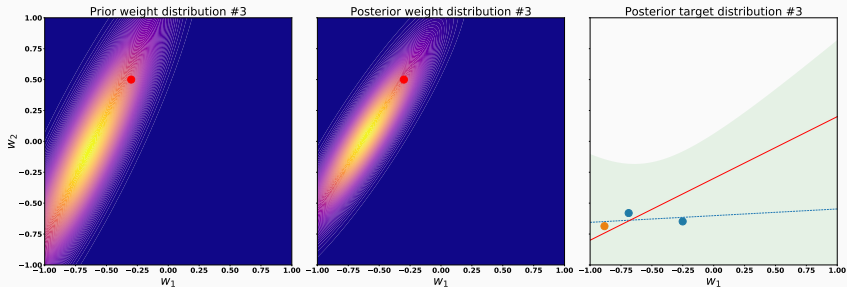
Online linear regression - progressive validation

$$y = w_1 + w_2x, \quad t = -0.3 + 0.5x + \epsilon, \quad \beta = 25, \alpha = 2$$



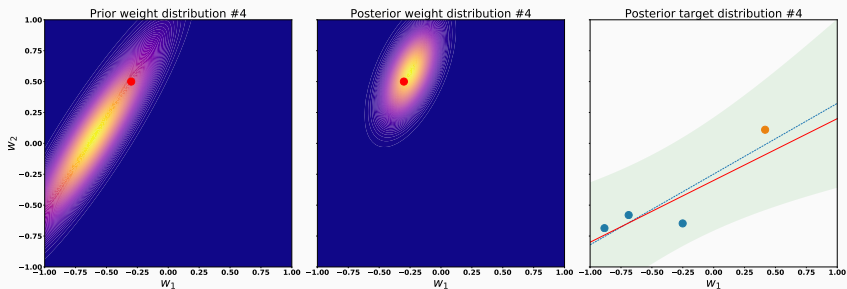
Online linear regression - progressive validation

$$y = w_1 + w_2x, \quad t = -0.3 + 0.5x + \epsilon, \quad \beta = 25, \alpha = 2$$



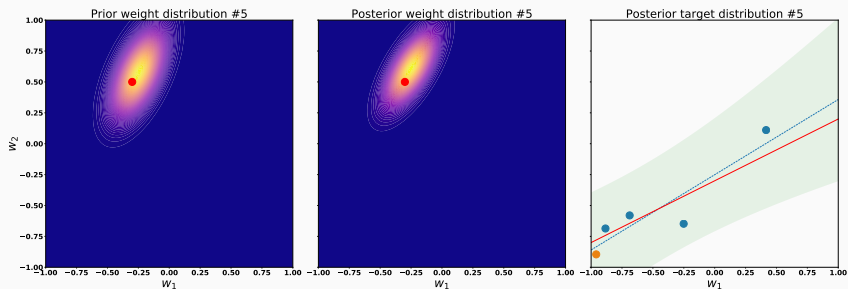
Online linear regression - progressive validation

$$y = w_1 + w_2x, \quad t = -0.3 + 0.5x + \epsilon, \quad \beta = 25, \alpha = 2$$



Online linear regression - progressive validation

$$y = w_1 + w_2x, \quad t = -0.3 + 0.5x + \epsilon, \quad \beta = 25, \alpha = 2$$





C. M. Bishop.

Pattern Recognition and Machine Learning.

Springer, 2006.