

UNIT 2

Bayes Techniques, pPCA



Johannes Brandstetter

Institute for Machine Learning

Copyright statement:

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Likelihood

- Likelihood is defined as:

$$\mathcal{L}(\{\mathbf{z}\}; \mathbf{w}) = p(\{\mathbf{z}\}; \mathbf{w}) = \prod_{i=1}^l p(\mathbf{z}_i; \mathbf{w})$$

- I.e., probability of the (data-distribution) model $p(\mathbf{z}; \mathbf{w})$ to (re-)produce the dataset
- In supervised learning we can write:

$$p(\mathbf{z}; \mathbf{w}) = p(\mathbf{x})p(y|\mathbf{x}; \mathbf{w})$$

- and

$$p(\{\mathbf{z}\}; \mathbf{w}) = \prod_{i=1}^l p(\mathbf{x}_i) \prod_{i=1}^l p(y_i|\mathbf{x}_i; \mathbf{w})$$

Likelihood

- Because $\prod_{i=1}^l p(\mathbf{x}_i)$ is independent of the parameters, it is sufficient to maximize the conditional likelihood:

$$p(y|\{\mathbf{x}\}; \mathbf{w}) = \prod_{i=1}^l p(y_i|\mathbf{x}_i; \mathbf{w})$$

- 2 views: (i) parameter vector \mathbf{w} is used to parameterize the likelihood, and (ii) **likelihood is conditioned on \mathbf{w}** !
- **Supervised machine learning**: We choose parameters \mathbf{w} such that they maximize the probability of the observed dataset $\{z_1, \dots, z_l\}$ given our choice of model \mathcal{M} :

$$p(\{z\}|\mathbf{w}, \mathcal{M}) = \prod_{i=1}^l p(z_i|\mathbf{w}, \mathcal{M}) .$$

- Usually \mathcal{M} is left out for simplicity

Frequentist approach

- From a probabilistic perspective, frequentists are trying to maximize the likelihood $p(\{z\}|\mathbf{w}, \mathcal{M})$
- A (statistical) model is simply a probability distribution over data \Rightarrow **Maximum likelihood estimation** (MLE) to obtain, or “train”, predictive models
- If the model is complex enough, only maximizing the likelihood **leads to overfitting**
- In the most extreme case the model produces the training examples with equal probability and other data with probability zero:
 $\Rightarrow p(\{z\}|\mathbf{w}, \mathcal{M})$ is the sum of Dirac delta-distributions

From Frequentist to the Bayesian approach

- E.g., to avoid overfitting we can assume that **certain w are more probable to be observed** in the real world than others
- That means some models are more likely
- The fact that some models are more likely can be expressed by a distribution $p(w)$, the **prior distribution**
- The information in $p(w)$ stems from prior knowledge about the problem

Bayesianists

- Crucial property of the Bayesian approach is to realistically **quantify uncertainty**
- Instead of a parameter point estimate, a Bayesian approach defines a full **probability distribution** over parameters \Rightarrow **posterior distribution**
- Posterior represents our **belief/hypothesis/uncertainty** about the value of each parameter (setting)

Bayes' theorem

- Bayes' theorem:

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

- Bayes' theorem for data modeling:

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}, \text{ or } p(\mathbf{w}|\{\mathbf{z}\}) = \frac{p(\{\mathbf{z}\}|\mathbf{w}) \times p(\mathbf{w})}{p_{\mathbf{w}}(\{\mathbf{z}\})}$$

- **Prior distribution** $p(\mathbf{w})$ over the parameters to capture our belief about what our model parameters should look like prior to observing any data
- Using our dataset, we can **update** (multiply) our prior belief with the **likelihood** $p(\{\mathbf{z}\}|\mathbf{w})$

Bayes' theorem cont'd

- Likelihood $p(\{z\}|w)$ is the same quantity we saw in the frequentist approach. It tells us how well the observed data is explained by a specific parameter setting w
- The product $p(\{z\}|w)p(w)$ must be evaluated for **each parameter** setting, and normalized:
⇒ **Marginalizing** (summing or integrating) over all parameter settings
- The normalization constant

$$p_w(\{z\}) = \int_W p(\{z\}|w)p(w)dw$$

is called **evidence** for a class of models

- $p_w(\{z\})$ provides evidence for how good our model as a whole, i.e., how likely the data is, taking into account all **parameter settings**

Example: Conjugate prior

- **Conjugate distribution** (conjugate pair): pair of (i) sampling distribution (likelihood) and (ii) prior distribution, for which the resulting posterior distribution belongs into the same parametric family of distributions than the prior distribution
- We say that the **prior distribution** is a conjugate prior for this sampling distribution
- Example: Gamma distribution is a conjugate prior of the Poisson distribution
- Conjugate priors give convenience: we can do statistical inference, e.g., hypothesis testing, with the posterior distribution

Conjugate pair: Gamma & Poisson distr

- Gamma distribution:

$$f(\mathbf{x}; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \mathbf{x}^{\alpha-1} e^{-\beta \mathbf{x}}, \quad \Gamma(n) = (n-1)!$$

- Likelihood (Poisson distribution):

$$p(\{\mathbf{x}\}|\lambda) = \prod_{i=1}^l p(\mathbf{x}_i|\lambda) = \prod_{i=1}^l \lambda^{\mathbf{x}_i} \frac{e^{-\lambda}}{\mathbf{x}_i!} \propto \lambda^{\sum_{i=1}^l \mathbf{x}_i} e^{-l\lambda}, \quad \mathbf{x}_i = \{0, 1, \dots\}$$

- We treat the term $\int \propto p(\{\mathbf{x}\}|\lambda)p(\lambda)$ as constant since it does not depend on λ

- **Unnormalized posterior distribution** for λ :

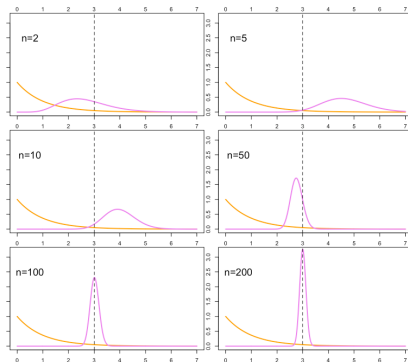
$$\begin{aligned} p(\lambda|\{\mathbf{x}\}) &\propto p(\{\mathbf{x}\}|\lambda)p(\lambda) \\ &\propto \lambda^{l\bar{\mathbf{x}}} e^{-l\lambda} \lambda^{\alpha-1} e^{-\beta\lambda} \\ &\propto \lambda^{\alpha+l\bar{\mathbf{x}}-1} e^{-(\beta+l)\lambda} \quad \text{with } \bar{\mathbf{x}} = \sum_{i=1}^l \mathbf{x}_i \end{aligned}$$

- The posterior distribution is thus:

$$p(\lambda|\{\mathbf{x}\}) \propto \Gamma(\alpha + l\bar{\mathbf{x}}, \beta + l)$$

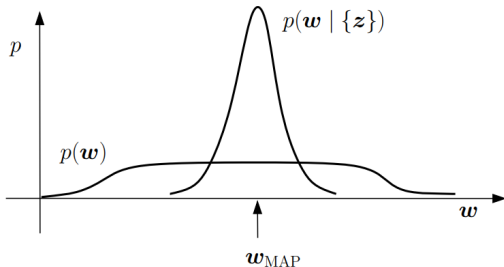
Conjugate pair: Gamma & Poisson distribution

- Example: Generate observations from Poisson distribution with a known parameter $\lambda = 3$
- Prior distribution: $\lambda \sim \text{Gamma}(\alpha = 1, \beta = 1)$
- Prior distribution: $\lambda \sim \text{Gamma}(\alpha = 1, \beta = 1)$



Example: Maximum A posteriori (MAP)

- The Maximum A Posteriori Approach (MAP) searches for the maximal posterior $p(\mathbf{w}|\{\mathbf{z}\})$



- MAP estimate picks the parameter setting that has the **highest probability** assigned to it (the distribution's mode)
- MAP provides a point estimate (not fully Bayesian)

MAP for Gaussian weight prior

- Gaussian weight prior:

$$p(\mathbf{w}) = \frac{1}{Z_{\mathbf{w}}(\alpha)} \exp\left(-\frac{1}{2}\alpha\|\mathbf{w}\|^2\right)$$
$$Z_{\mathbf{w}}(\alpha) = \int_{\mathbf{w}} \exp\left(-\frac{1}{2}\alpha\|\mathbf{w}\|^2\right) d\mathbf{w} = \left(\frac{2\pi}{\alpha}\right)^{\frac{W}{2}}$$

- The negative log-posterior is:

$$-\ln p(\mathbf{w}|\{\mathbf{z}\}) = -\ln(\{\mathbf{z}\}|\mathbf{w}) - \ln p(\mathbf{w}) + \ln p_{\mathbf{w}}(\{\mathbf{z}\})$$

- For MAP estimation only $-\ln(\{\mathbf{z}\}|\mathbf{w}) - \ln p(\mathbf{w})$ must be minimized which results in the terms:

$$R(\mathbf{w}) = \underbrace{\beta R_{\text{emp}}}_{\text{empirical error}} + \underbrace{\alpha \Omega(\mathbf{w})}_{\text{complexity term}}$$

- The complexity term $\Omega(\mathbf{w}) = \|\mathbf{w}\|^2$ can be viewed in most cases as MAP estimation

MAP for Gaussian weight prior

- For Gaussian noise models, the likelihood is the exponential function with empirical error as an argument:

$$p(\{z|\mathbf{w}\}) = \frac{1}{Z_R(\beta)} \exp\left(-\frac{1}{2}\beta R_{\text{emp}}\right) ,$$

- the prior is an exponential function of the complexity:

$$p(\mathbf{w}) = \frac{1}{Z_W(\alpha)} \exp\left(-\frac{1}{2}\alpha\Omega(\mathbf{w})\right) ,$$

- and the posterior is:

$$p(\mathbf{w}|\{z\}) = \frac{1}{Z(\alpha, \beta)} \exp\left(-\frac{1}{2}(\alpha\Omega(\mathbf{w}) + \beta R_{\text{emp}})\right)$$
$$Z(\alpha, \beta) = \int_W \exp\left(-\frac{1}{2}(\alpha\Omega(\mathbf{w}) + \beta R_{\text{emp}})\right) d\mathbf{w}$$

Evidence

- The full fledged Bayesian approach is to specify a **predictive distribution**, i.e., the evidence:

$$p(y|\{z\}, x) = \int p(y|w, x)p(w|\{z\})dw$$

- This defines the probability for class label y given new input x and dataset $\{z\}$
- To compute the evidence we need to **marginalize over all parameter settings**
- We multiply the posterior probability of each setting w with the probability of label y given input x using parameter setting $w \Rightarrow$ **Bayesian Model Averaging (BMA)**

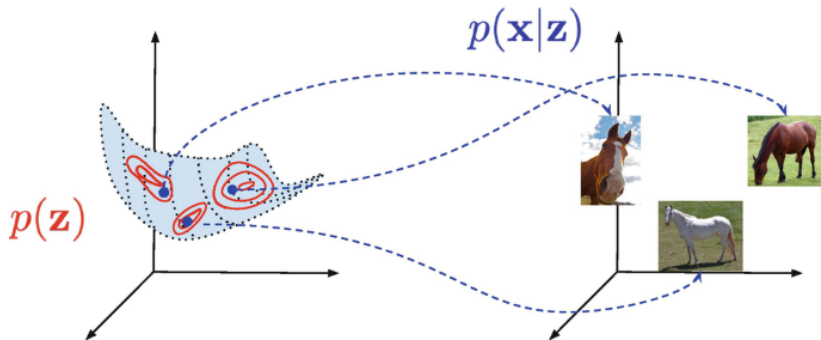
Marginalization

- In practice marginalizing over the whole parameter space is often impossible (**intractable**) as we can have infinitely many such settings \Rightarrow Bayesian approach is fundamentally about marginalization
- Different family of methods to learn parameter values:
 - **Sampling methods** such as Markov Chain Monte Carlo (MCMC)
 - **Variational inference** techniques use a simpler, tractable family of distributions. Prominent examples are **Variational Autoencoders**

Probabilistic Principal Component Analysis

Latent variable models

Deep Generative Modeling, Tomczak



Latent variable models

- Describe generative process as:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

- How to calculate this integral?
 - Tractable: Integrate
 - Intractable: Variational inference
- !! Unfortunately, literature uses \mathbf{z} for latent variables. This coincides with $\{\mathbf{z}\}$ which we used so far to describe dataset !!

Principal Component Analysis

- Mean and covariance of dataset $\{x_1, \dots, x_l\}$, where $x_i = (x_i^{(1)}, \dots, x_i^{(d)})^T \in \mathbb{R}^d$:

$$\bar{x} = \frac{1}{l} \sum_{i=1}^l x_i, \quad S = \frac{1}{l} \sum_{i=1}^l (x_i - \bar{x})(x_i - \bar{x})^T$$

- We consider a unit vector $u_1 \in \mathbb{R}^d$, i.e., $u_1^T u_1 = 1$
- $u_1^T x_i$ projects each data point onto a scalar value
- The variance of the projected data is given by:

$$\frac{1}{l} \sum_{i=1}^l (u_1^T x_i - u_1^T \bar{x})^2 = u_1^T S u_1$$

Principal Component Analysis

- We **maximize the projected variance** $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ wrt \mathbf{u}_1
- Normalization condition $\mathbf{u}_1^T \mathbf{u}_1 = 1 \Rightarrow$ we enforce normalization constraint with Lagrangian multiplier α_1
- Unconstrained maximization of

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \alpha_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

- Stationary point when:

$$\mathbf{S} \mathbf{u}_1 = \alpha_1 \mathbf{u}_1$$

- $\Rightarrow \mathbf{u}_1$ must be an eigenvector of \mathbf{S}
- Using $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \alpha_1$: projected variance is maximum when \mathbf{u}_1 equals the eigenvector of \mathbf{S} with the largest eigenvalue $\lambda_1 \Rightarrow \mathbf{u}_1$ is the first principal component of \mathbf{S}
- PCA: **linear projection** of the data onto a subspace of lower dimensionality than the original data space

Probabilistic Principal Component Analysis

- All of the marginal and conditional distributions are Gaussian:
 - **Latent variable** $z \in \mathbb{R}^m$ corresponding to the m -**dimensional principal-component** subspace
 - Gaussian prior distribution $p(z)$ over the latent variable
 - Gaussian conditional distribution $p(x|z) = \mathcal{N}(\mathbf{W}z + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ for the observed data $\{x_1, \dots, x_l\} \in \mathbb{R}^d$ conditioned on z :

$$p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

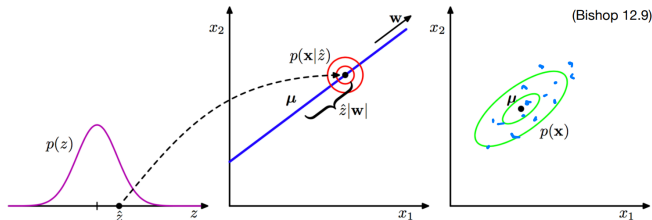
$$p(x|z) = \mathcal{N}(\mathbf{W}z + \boldsymbol{\mu}, \sigma^2 \mathbf{I}), \quad \mathbf{W} \in \mathbb{R}^{d \times m}, \quad \boldsymbol{\mu} \in \mathbb{R}^d, \quad \mathbf{I} \in \mathbb{R}^{d \times d}$$

- ... assuming a linear dependency between z and x :

$$x = \mathbf{W}z + \boldsymbol{\mu} + \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon|\mathbf{0}, \sigma^2 \mathbf{I})$$

- The columns of \mathbf{W} span a **linear subspace** within the data space that corresponds to the principal subspace

Probabilistic Principal Component Analysis



- We choose a value for the latent z
- ... and then sample x from an isotropic Gaussian distribution with mean $\mathbf{W}z + \mu$ and covariance $\sigma^2\mathbf{I}$

Probabilistic Principal Component Analysis

- The marginal distribution $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ is also Gaussian (sum of Gaussians is Gaussian, product of Gaussians is Gaussian), i.e., $\mathcal{N}(\mathbf{x}|\cdot, \cdot)$

$$\begin{aligned} E_{p(\mathbf{x}|\mathbf{z})}(\mathbf{x}) &= E_{p(\mathbf{x}|\mathbf{z})}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}) \\ &= \mathbf{W} \underbrace{E_{p(\mathbf{x}|\mathbf{z})}(\mathbf{z})}_{=0} + E_{p(\mathbf{x}|\mathbf{z})}(\boldsymbol{\mu}) + \underbrace{E_{p(\mathbf{x}|\mathbf{z})}(\boldsymbol{\epsilon})}_{=0} \\ &= \boldsymbol{\mu} \end{aligned}$$

$$\begin{aligned} \text{Cov}_{p(\mathbf{x}|\mathbf{z})}(\mathbf{x}) &= E_{p(\mathbf{x}|\mathbf{z})} \left((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right) \\ &= E_{p(\mathbf{x}|\mathbf{z})} \left((\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} - \boldsymbol{\mu})(\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon} - \boldsymbol{\mu})^T \right) \\ &= E_{p(\mathbf{x}|\mathbf{z})} \left(\mathbf{W}\mathbf{z}\mathbf{z}^T\mathbf{W}^T + \mathbf{W}\mathbf{z}\boldsymbol{\epsilon}^T + \boldsymbol{\epsilon}\mathbf{z}^T\mathbf{W}^T + \boldsymbol{\epsilon}\boldsymbol{\epsilon}^T \right) \\ &= \underbrace{\mathbf{W} E(\mathbf{z}\mathbf{z}^T) \mathbf{W}^T}_{=\mathbf{I}} + \underbrace{\mathbf{W} E(\mathbf{z}\boldsymbol{\epsilon}^T)}_{=\mathbf{W}E(\mathbf{z})E(\boldsymbol{\epsilon}^T)} + \underbrace{E(\boldsymbol{\epsilon}\mathbf{z}^T)\mathbf{W}^T}_{=E(\boldsymbol{\epsilon})E(\mathbf{z}^T)\mathbf{W}^T} + \underbrace{E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T)}_{=\sigma^2\mathbf{I}} \end{aligned}$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$$

True posterior

- Due to properties of Gaussians (Bishop 2.116), we can also calculate the true posterior $p(z|x)$

$$p(z|x) = \mathcal{N}\left((\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^T (\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2} (\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I})\right)$$

- If we find \mathbf{W} , we can calculate $p(z|x) \Rightarrow$ maximum-likelihood estimation
- For a given observation x , we can calculate the **distribution over the latent factors**

Maximum likelihood PCA

- Log-likelihood of $p(\mathbf{x})$:

$$\begin{aligned}\ln(p(\mathbf{x})) &= \sum_{i=1}^l \ln \left(p(\mathbf{x}_i | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}) \right) \\ &= -\frac{ld}{2} \ln(2\pi) - \frac{l}{2} \ln \left(|\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}| \right) \\ &\quad - \frac{1}{2} \sum_{i=1}^l (\mathbf{x} - \boldsymbol{\mu})^T \left(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I} \right)^{-1} (\mathbf{x} - \boldsymbol{\mu})\end{aligned}$$

- Once again, we take the derivative w.r.t. the parameter of interest, set it to zero, and solve it. For the mean:

$$\begin{aligned}\sum_{i=1}^l \left(\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I} \right)^{-1} (\mathbf{x} - \boldsymbol{\mu}) &= 0 \\ \Rightarrow \boldsymbol{\mu} &= \frac{1}{l} \sum_{i=1}^l \mathbf{x}_i = \bar{\mathbf{x}}\end{aligned}$$

Maximum likelihood PCA

- Maximization with respect to \mathbf{W} and σ^2 is more complex, but nonetheless has an **exact closed-form solution** (Tipping & Bishop 1999)

$$\mathbf{W}_{mk} = \mathbf{U}_m (\mathbf{L}_m - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}$$

- \mathbf{U}_m is a $d \times m$ matrix whose columns are given by any subset (of size M) of the eigenvectors of the data covariance matrix \mathbf{S}
- \mathbf{L}_m is an $m \times m$ diagonal matrix with elements corresponding to the eigenvalues
- \mathbf{R} is an arbitrary $m \times m$ orthogonal matrix
- **Maximum of the likelihood function is obtained when the m eigenvectors are chosen to be those whose eigenvalues are the m largest**

Probabilistic PCA has closed-form solutions

- Probabilistic PCA is the probabilistic generative version of PCA: we can also draw samples from it
- Probabilistic PCA is a form of Gaussian distribution with number of parameters restricted (by the latent space)
- Limitations:
 - Latent space $p(z) = \mathcal{N}(z|\mathbf{0}, \mathbf{I})$ is very simple
 - For generation of $p(x|z)$ only **linear transformations** possible