

# Learning with Gaussians

CPSC 440/550: Advanced Machine Learning

`cs.ubc.ca/~dsuth/440/23w2`

University of British Columbia, on unceded Musqueam land

2023-24 Winter Term 2 (Jan–Apr 2024)

- Project guidelines finally available
- Brief proposal due March 29th
  - 10% of project grade, “lightly graded”: mostly checking scope of project
  - If you hand in earlier, we’ll give you scope feedback earlier
- Actual project due last day of finals (Saturday, April 27)
  - 6-page writeup, plus possible appendices/code supplement
  - Details on format to come

## Last time: Multivariate Gaussians

- Continuous density estimation,  $d > 1$  with the multivariate Gaussian distribution

$$x \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \text{means} \quad p(x \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(x - \boldsymbol{\mu})\right)$$

- If  $\boldsymbol{\Sigma}$  is a diagonal matrix, product of univariate normals
- $\boldsymbol{\mu}_j$  is  $\mathbb{E}[x_j]$ ;  $\boldsymbol{\Sigma}_{jj'}$  gives  $\text{Cov}(x_j, x_{j'})$ 
  - If  $\text{Cov}(x_j, x_{j'}) = 0$ , then  $x_j \perp x_{j'}$  (for jointly-Gaussian variables)
- If  $\boldsymbol{\Sigma}$  is singular, “degenerate” Gaussian:  $v^T x$  takes a constant value for some  $v$
- $Ax + b \sim \mathcal{N}(A\boldsymbol{\mu} + b, A\boldsymbol{\Sigma}A^T)$ 
  - Lets us sample based on  $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - Marginalizing: still normal, just ignore the other variables in  $\boldsymbol{\mu}, \boldsymbol{\Sigma}$
  - Conditioning:  $x \mid z \sim \mathcal{N}(\boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xz}\boldsymbol{\Sigma}_z^{-1}(z - \boldsymbol{\mu}_z), \boldsymbol{\Sigma}_x - \boldsymbol{\Sigma}_{xz}\boldsymbol{\Sigma}_z^{-1}\boldsymbol{\Sigma}_{xz}^T)$

# Outline

- 1 Learning multivariate Gaussians
- 2 Generative classifiers with Gaussians
- 3 Bayesian Linear Regression

## MLE for the mean of a multivariate Gaussian

- If  $x^{(i)} \stackrel{iid}{\sim} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  for  $\boldsymbol{\Sigma} \succ 0$ , we have

$$p\left(x^{(i)} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \left(x^{(i)} - \boldsymbol{\mu}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(x^{(i)} - \boldsymbol{\mu}\right)\right),$$

so up to a constant our **negative log-likelihood** for  $n$  examples is

$$\frac{1}{2} \sum_{i=1}^n \left(x^{(i)} - \boldsymbol{\mu}\right)^{\top} \boldsymbol{\Sigma}^{-1} \left(x^{(i)} - \boldsymbol{\mu}\right) + \frac{n}{2} \log |\boldsymbol{\Sigma}|$$

- This is a **convex quadratic** in  $\boldsymbol{\mu}$ ; setting gradient to zero gives

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

- Mean along each dimension; it **doesn't depend on  $\boldsymbol{\Sigma}$**

## MLE for the covariance of a multivariate Gaussian

- To get MLE for  $\Sigma$  we can re-parameterize in terms of **precision matrix**  $\Theta = \Sigma^{-1}$ ,

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^n \left( x^{(i)} - \hat{\mu} \right)^{\top} \Sigma^{-1} \left( x^{(i)} - \hat{\mu} \right) + \frac{n}{2} \log |\Sigma| \\ = \frac{1}{2} \sum_{i=1}^n \left( x^{(i)} - \hat{\mu} \right)^{\top} \Theta \left( x^{(i)} - \hat{\mu} \right) + \frac{n}{2} \log |\Theta^{-1}| \end{aligned}$$

- After some work (**bonus slides**), we get that this is equal to

$$f(\Theta) = \frac{n}{2} \text{Tr}(\mathbf{S}\Theta) - \frac{n}{2} \log |\Theta|, \text{ with } \mathbf{S} = \frac{1}{n} \sum_{i=1}^n \left( x^{(i)} - \hat{\mu} \right) \left( x^{(i)} - \hat{\mu} \right)^{\top}$$

- $\mathbf{S}$  is the **sample covariance**: if  $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}_n \hat{\mu}^{\top}$  is centred data,  $S = \frac{1}{n} \tilde{\mathbf{X}}^{\top} \tilde{\mathbf{X}}$
- **Trace operator**  $\text{Tr}(A)$  is the sum of the diagonal elements of  $A$
- $\text{Tr}(A^{\top} B) = \sum_j (A^{\top} B)_{jj} = \sum_j \sum_i (A^{\top})_{ji} B_{ij} = \sum_{ij} A_{ij} B_{ij}$ , i.e.  $(A * B) . \text{sum}()$

## MLE for the covariance of a multivariate Gaussian

- Gradient matrix of NLL with respect to  $\Theta$  is (not obvious, see **bonus slides**)

$$\nabla f(\Theta) = \frac{n}{2} (\mathbf{S} - \Theta^{-1}) \quad \text{for } S = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \hat{\mu})(x^{(i)} - \hat{\mu})^T$$

- The MLE for a given  $\mu$  is obtained by setting the gradient matrix to zero, giving

$$\Theta = \mathbf{S}^{-1} \quad \text{or} \quad \Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \hat{\mu})(x^{(i)} - \hat{\mu})^T$$

- To have  $\Sigma \succ 0$ , we **need a positive-definite sample covariance,  $S \succ 0$** 
  - If  $S$  is not positive definite, NLL is unbounded below, and MLE doesn't exist
  - Like requiring “not all values are the same” in univariate Gaussian
  - In  $d$  dimensions, you need  $d$  linearly independent  $x^{(i)}$  values (no “multi-collinearity”)
  - This is only possible if  $n \geq d$ ! (But might not be true even if it is)
- Note: most distributions' MLEs **don't** correspond with “moment matching”

## Example: Multivariate Gaussians on MNIST

- Let's try **continuous** density estimation on (binary) handwritten digits

$$x^i = \text{vec} \left( \begin{array}{c} \text{[Handwritten digit 4]} \end{array} \right) \in \mathbb{R}^{784}$$

Diagonal  $\Sigma$ :

$$\hat{\mu} = \text{vec} \left( \begin{array}{c} \text{[Blurred digit 4]} \end{array} \right)$$

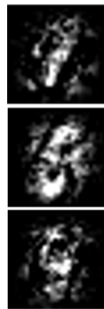
$$\hat{\Sigma} = \text{diag} \left( \text{vec} \left( \begin{array}{c} \text{[Blurred digit 4]} \end{array} \right) \right)$$



General  $\Sigma$ :

$\hat{\mu}$  is the same (!)  
 $\hat{\Sigma}$  is big  
(784 by 784)

7 7  
4 4



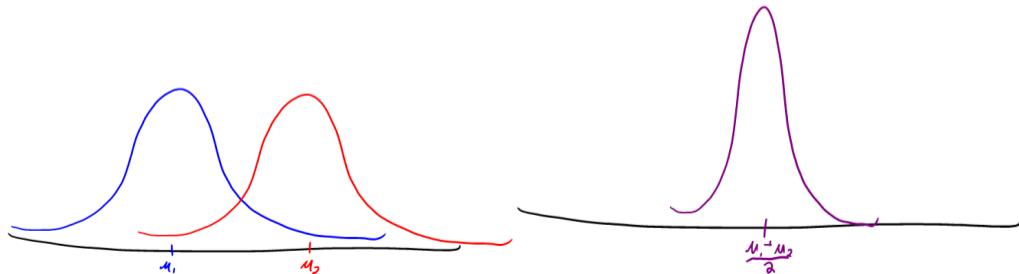


## Product of Gaussian densities

- This property will be helpful in deriving MAP/Bayesian estimation:
- Consider a variable  $x$  whose pdf is written as product of two Gaussians,

$$p(x) \propto \underbrace{\mathcal{N}(x \mid \boldsymbol{\mu}_1, \mathbf{I})}_{\text{density of } \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{I}) \text{ at } x} \mathcal{N}(x \mid \boldsymbol{\mu}_2, \mathbf{I})$$

- This **product of Gaussian pdfs is a Gaussian** with  $\boldsymbol{\mu} = \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)$  and  $\boldsymbol{\Sigma} = \frac{1}{2}\mathbf{I}$



## Product of Gaussian densities

- If  $p(x) \propto \mathcal{N}(x | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \mathcal{N}(x | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ ,  
then  $x$  is Gaussian with (see PML2 2.2.7.6 – complete the square in the exponent)

$$\text{covariance } \boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}$$

$$\text{mean } \boldsymbol{\mu} = \boldsymbol{\Sigma} \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma} \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2$$

- Consider  $x^{(i)} \sim \mathcal{N}(x^{(i)} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  for fixed  $\boldsymbol{\Sigma}$  and  $\boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ :

$$p(\boldsymbol{\mu} | \mathbf{X}, \boldsymbol{\Sigma}, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \propto p(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \prod_{i=1}^n p(x^{(i)} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (\text{Bayes rule})$$

$$= p(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \prod_{i=1}^n p(\boldsymbol{\mu} | x^{(i)}, \boldsymbol{\Sigma}) \quad (\text{symmetry of } x^{(i)} \text{ and } \boldsymbol{\mu})$$

$$= (\text{product of } (n + 1) \text{ Gaussians})$$

- So, working it out gives. . .

## MAP estimation for mean

- For fixed  $\Sigma$ , conjugate prior for mean is a Gaussian:

$$x^{(i)} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma) \quad \boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0) \quad \text{implies} \quad \boldsymbol{\mu} \mid \mathbf{X}, \Sigma \sim \mathcal{N}(\boldsymbol{\mu}^+, \Sigma^+),$$

where

$$\Sigma^+ = (n\Sigma^{-1} + \Sigma_0^{-1})^{-1},$$

$$\boldsymbol{\mu}^+ = \Sigma^+ (n\Sigma^{-1} \boldsymbol{\mu}_{\text{MLE}} + \Sigma_0^{-1} \boldsymbol{\mu}_0) \quad \text{MAP estimate of } \boldsymbol{\mu}$$

- In special case of  $\Sigma = \sigma^2 \mathbf{I}$  and  $\Sigma_0 = \frac{1}{\lambda} \mathbf{I}$ , we get

$$\Sigma^+ = \left( \frac{n}{\sigma^2} \mathbf{I} + \lambda \mathbf{I} \right)^{-1} = \frac{1}{\frac{n}{\sigma^2} + \lambda} \mathbf{I},$$

$$\boldsymbol{\mu}^+ = \Sigma^+ \left( \frac{n}{\sigma^2} \boldsymbol{\mu}_{\text{MLE}} + \lambda \boldsymbol{\mu}_0 \right)$$

- Posterior predictive is  $\mathcal{N}(\boldsymbol{\mu}^+, \Sigma + \Sigma^+)$  – take product of  $(n+2)$  then marginalize
  - Many Bayesian inference tasks have closed form; if not, Monte Carlo is easy

## MAP Estimation in Multivariate Gaussian (Trace Regularization)

- A common MAP estimate for  $\Sigma$  is

$$\hat{\Sigma} = \mathbf{S} + \lambda \mathbf{I},$$

where  $S$  is the covariance of the data

- Key advantage:  $\hat{\Sigma}$  is strictly positive definite (eigenvalues are at least  $\lambda$ )
- This corresponds to L1 regularization of precision diagonals (see bonus)

$$f(\Theta) = \underbrace{\text{Tr}(\mathbf{S}\Theta) - \log |\Theta|}_{\text{NLL times } 2/n} + \lambda \sum_{j=1}^d |\Theta_{jj}|$$

- Note this *doesn't* set  $\Theta_{jj}$  values to exactly zero
  - Log-determinant term becomes arbitrarily steep as the  $\Theta_{jj}$  approach 0
  - It's not really the case that "L1 gives sparsity"; it's "L2 + L1 gives sparsity"

# Conjugate Priors for Covariance

bonus!

- Trace regularization (or Graphical LASSO, later): **not a conjugate prior**
- Conjugate prior for  $\Theta$  with known mean is **Wishart** distribution
  - A multi-dimensional **generalization of the gamma** distribution
    - Gamma is a distribution over positive scalars
    - Wishart is a **distribution over positive-definite matrices**
  - Posterior predictive is a **student  $t$**  distribution
  - Conjugate prior for  $\Sigma$  is **inverse-Wishart** (equivalent posterior)
- If both  $\mu$  and  $\Theta$  are variables, conjugate prior is **normal-Wishart**
  - Normal times Wishart, with a particular dependency among parameters
  - Posterior predictive is again a **student  $t$**  distribution
- Wikipedia has already done a lot of possible homework questions for you:
  - [https://en.wikipedia.org/wiki/Conjugate\\_prior](https://en.wikipedia.org/wiki/Conjugate_prior)

# Outline

- 1 Learning multivariate Gaussians
- 2 Generative classifiers with Gaussians**
- 3 Bayesian Linear Regression

# Generative Classification with Gaussians

- Consider a **generative classifier** with **continuous features**:

$$p(y | x) \propto p(x, y) = \underbrace{p(x | y)}_{\text{continuous}} \underbrace{p(y)}_{\text{discrete}}$$

- Model  $y$  as a categorical distribution (classification task)
- Previously handled  $p(x | y)$  with the **naive Bayes** assumption,  $x_i \perp\!\!\!\perp x_j | y$ 
  - Strong, usually unrealistic assumption
- In **Gaussian discriminant analysis (GDA)** we assume  $x | y$  is **Gaussian**
  - Classifier asks “which Gaussian makes this  $x^{(i)}$  most likely?”
  - This can **model pairwise correlations** within each class
    - Doesn't need the naive Bayes assumption

# Gaussian Discriminant Analysis (GDA)

- In **Gaussian discriminant analysis** we assume  $x | y$  is Gaussian

$$p(x, y = c) = \underbrace{p(y) p(x | y = c)}_{\text{product rule}} = \underbrace{\pi_c}_{\text{Pr}(y=c)} \underbrace{p(x | \mu_c, \Sigma_c)}_{\text{Gaussian pdf}}$$

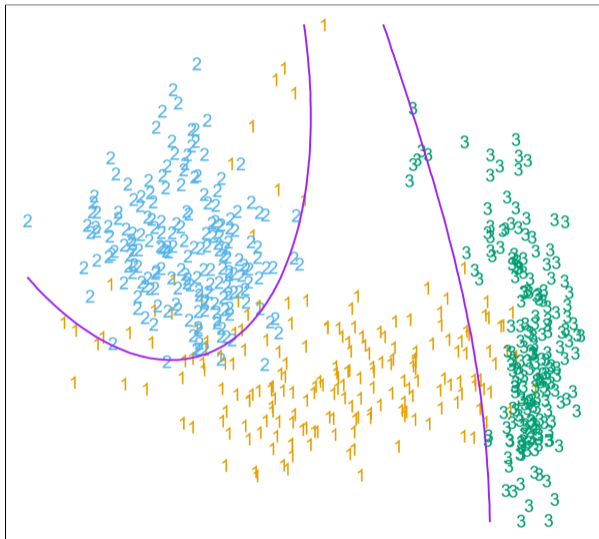
- Classify based on

$$\begin{aligned} \arg \max_c p(y = c | x) &= \arg \max_c \log p(y = c, x) \\ &= \arg \max_c \log \pi_c - \frac{1}{2} \log |\Sigma_c| - \frac{1}{2} (x - \mu_c)^\top \Sigma_c^{-1} (x - \mu_c) \end{aligned}$$

- With general choices for  $\mu_c$  and  $\Sigma_c$ , we're taking the **max of  $k$  quadratics**
  - Means that the decision boundary will be zeros of a quadratic ("quadric surface")
  - Leads to the equivalent name **quadratic discriminant analysis** (QDA)
- Fitting GDA=QDA: fit  $\pi_c$  as categorical, fit Gaussian for each subset with  $y^{(i)} = c$



# GDA=QDA example



## Special case: Linear Discriminant Analysis (LDA)

- A common special case: constrain  $\Sigma_c = \Sigma$  for all  $c$
- Means that we classify as

$$\begin{aligned}\arg \max_c p(y = c | x) &= \arg \max_c \log \pi_c - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x - \mu_c)^\top \Sigma^{-1} (x - \mu_c) \\ &= \arg \max_c \log \pi_c - \frac{1}{2} x^\top \Sigma^{-1} x + \mu_c^\top \Sigma^{-1} x - \frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c \\ &= \arg \max_c \underbrace{(\Sigma^{-1} \mu_c)^\top x}_{w_c} + \underbrace{\log \pi_c - \frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c}_{b_c}\end{aligned}$$

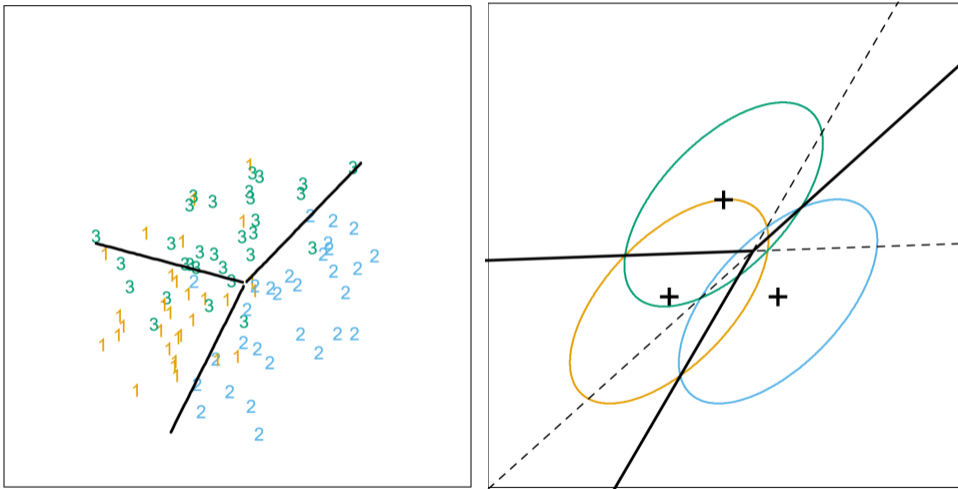
so this is a linear classifier!

- Behaves (asymptotically) optimally **if the assumptions are true**:  $x | y \sim \mathcal{N}(\mu_y, \Sigma)$
- May be terrible if these assumptions aren't true
- MLE in this model is simple:  $\mu_c$  is mean of the points with  $y^{(i)} = c$ ,

$$\Sigma \text{ is } \frac{1}{n} \sum_{i=1}^n \left( x^{(i)} - \mu_{y^{(i)}} \right) \left( x^{(i)} - \mu_{y^{(i)}} \right)^\top$$

## LDA example

- Example of fitting linear discriminant analysis (LDA) to a 3-class problem:



- We classify according to

$$\begin{aligned} \arg \max_c (\boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\mu}_c)^\top (\boldsymbol{\Sigma}^{-\frac{1}{2}} x) - \frac{1}{2} (\boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\mu}_c)^\top (\boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\mu}_c) + \log \pi_c \\ = \arg \max_c -\frac{1}{2} \|\boldsymbol{\Sigma}^{-\frac{1}{2}} x\|^2 + (\boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\mu}_c)^\top (\boldsymbol{\Sigma}^{-\frac{1}{2}} x) - \frac{1}{2} \|\boldsymbol{\Sigma}^{-\frac{1}{2}} \boldsymbol{\mu}_c\|^2 + \log \pi_c \\ = \arg \min_c \|\boldsymbol{\Sigma}^{-\frac{1}{2}} (x - \boldsymbol{\mu}_c)\|^2 - 2 \log \pi_c \end{aligned}$$

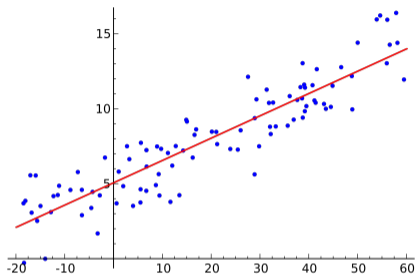
- If  $\pi_c$  are constant (all  $\frac{1}{k}$ ) and  $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$ , this picks the closest class mean
- With constant  $\pi_c$  but general  $\boldsymbol{\Sigma}$ , picks closest class mean in Mahalanobis distance

# Outline

- 1 Learning multivariate Gaussians
- 2 Generative classifiers with Gaussians
- 3 Bayesian Linear Regression**

# Regression with Gaussians

- In regression,  $y$  is continuous



[https://en.wikipedia.org/wiki/Regression\\_analysis](https://en.wikipedia.org/wiki/Regression_analysis)

- It's possible to use generative regression models (bonus slide)
  - For example, we could model  $p(x, y)$  as a multivariate Gaussian
    - Then use that the conditional  $p(y | x)$  is Gaussian for prediction
- But we usually treat features as fixed (as in discriminative classification models)
- Now ready to return to Bayesian linear regression

# Bayesian Linear Regression

- Linear regression with **Gaussian likelihood and prior**,

$$y | x \sim \mathcal{N}(w^\top x, \sigma^2), \quad w \sim \mathcal{N}(0, \lambda^{-1} \mathbf{I})$$

- MAP estimate is ridge regression (L2-regularized least squares)
- Can use Gaussian identities to work out that the **posterior** has the form

$$w | (\mathbf{X}, \mathbf{y}) \sim \mathcal{N} \left( w_{\text{MAP}}, \left( \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \right),$$

which is a **multivariate Gaussian centred at**  $w_{\text{MAP}} = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\sigma^2} \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{y}$

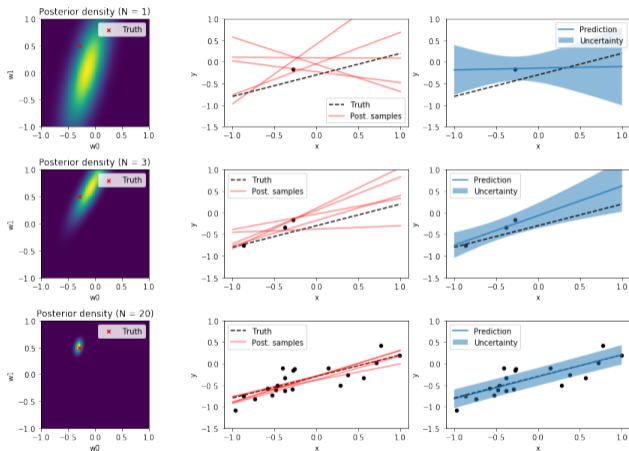
- The variance tells us **how much variation** we have around the MAP estimate
  - **In other models, the posterior mode (MAP) is usually not the posterior mean**
- By more Gaussian identities, the **posterior predictive** has the form

$$\tilde{y} | (\mathbf{X}, \mathbf{y}, \tilde{x}) \sim \mathcal{N} \left( w_{\text{MAP}}^\top \tilde{x}, \sigma^2 + \tilde{x}^\top \left( \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \tilde{x} \right)$$

- Posterior predictive mode=mean again the MAP prediction **in this model**
  - Working with the full posterior predictive gives us **variance of predictions**

# Bayesian Linear Regression

- Bayesian perspective gives us **variability in  $w$  and predictions**:

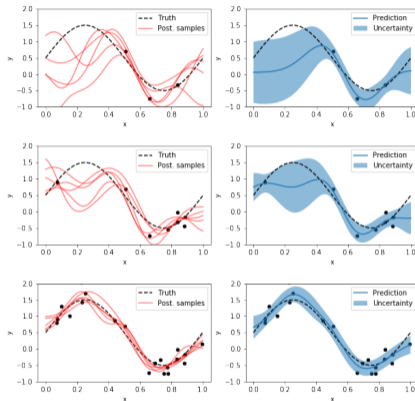


<http://krasserm.github.io/2019/02/23/bayesian-linear-regression>



# Bayesian Linear Regression

- Bayesian linear regression with Gaussian RBFs as features:



<http://krasserm.github.io/2019/02/23/bayesian-linear-regression>

- We have not only a prediction, but Bayesian inference gives “error bars”
  - Gives an idea of “where model is confident” and where it is not

## Digression: Gaussian Processes

bonus!

- In CPSC 340 you saw the **kernel trick**:
  - Rewrites L2-regularized least squares linear/prediction in terms of inner products
  - Allows us to efficiently use some exponential-sized or infinite-sized feature sets
- We can use **kernel trick on posterior** in Gaussian likelihood/prior model
  - Allows us to efficiently use some large or infinite-sized feature sets
  - Posterior in this case can be written as a **Gaussian process (GP)**
- Notation: a **stochastic process** is an **infinite collection of random variables**
- In a Gaussian process, **any finite subcollection is jointly Gaussian**
  - Defined in terms of a **mean function** and a **covariance function**
    - The set of **possible covariance functions is the set of possible kernel functions**
  - A popular book on this topic if you want to read more:  
Rasmussen/Williams, Gaussian Processes for Machine Learning
- We'll **assume we have explicit features**, but you could use kernels/GPs instead

# Summary

- Gaussian discriminant analysis and special case linear discriminant analysis
  - Generative classifier where  $x | y$  is multivariate normal
- Bayesian Linear Regression
  - Gaussian conditional likelihood and Gaussian prior gives Gaussian posterior
  - Posterior predictive is also Gaussian (“regression with error bars”)
- Next time: choosing priors, sampling from complex posteriors

- To get MLE for  $\Sigma$  we re-parameterize in terms of **precision matrix**  $\Theta = \Sigma^{-1}$ ,

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^n (x^{(i)} - \mu)^\top \Sigma^{-1} (x^i - \mu) + \frac{n}{2} \log |\Sigma| \\ &= \frac{1}{2} \sum_{i=1}^n (x^{(i)} - \mu)^\top \Theta (x^i - \mu) + \frac{n}{2} \log |\Theta^{-1}| \quad (\text{okay because } \Sigma \text{ is invertible}) \\ &= \frac{1}{2} \sum_{i=1}^n \text{Tr} \left( (x^{(i)} - \mu)^\top \Theta (x^i - \mu) \right) + \frac{n}{2} \log |\Theta|^{-1} \quad (\text{scalar } y^\top A y = \text{Tr}(y^\top A y)) \\ &= \frac{1}{2} \sum_{i=1}^n \text{Tr}((x^{(i)} - \mu)(x^i - \mu)^\top \Theta) - \frac{n}{2} \log |\Theta| \quad (\text{Tr}(ABC) = \text{Tr}(CAB)) \end{aligned}$$

- $|A^{-1}| = 1/|A|$  (can see e.g. from eigenvalues)
- The **trace** is the sum of the diagonal elements:  $\text{Tr}(A) = \sum_i A_{ii}$ 
  - $\text{Tr}(AB) = \text{Tr}(BA)$  when dimensions match: called **trace rotation** or **cyclic property**

- From the last slide,

$$p(\mathbf{X} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2} \sum_{i=1}^n \text{Tr} \left( \left( x^{(i)} - \hat{\boldsymbol{\mu}} \right) \left( x^{(i)} - \hat{\boldsymbol{\mu}} \right)^\top \boldsymbol{\Theta} \right) - \frac{n}{2} \log |\boldsymbol{\Theta}|$$

- We can **exchange the sum and trace** (trace is a linear operator) to get,

$$\begin{aligned} &= \frac{1}{2} \text{Tr} \left( \sum_{i=1}^n \left( x^{(i)} - \hat{\boldsymbol{\mu}} \right) \left( x^{(i)} - \hat{\boldsymbol{\mu}} \right)^\top \boldsymbol{\Theta} \right) - \frac{n}{2} \log |\boldsymbol{\Theta}| & \sum_i \text{Tr}(A_i B) &= \text{Tr} \left( \sum_i A_i B \right) \\ &= \frac{n}{2} \text{Tr} \left( \left( \underbrace{\left( \frac{1}{n} \sum_{i=1}^n \left( x^{(i)} - \hat{\boldsymbol{\mu}} \right) \left( x^{(i)} - \hat{\boldsymbol{\mu}} \right)^\top \right)}_{\text{sample covariance, } S} \right) \boldsymbol{\Theta} \right) - \frac{n}{2} \log |\boldsymbol{\Theta}| & \left( \sum_i A_i B \right) &= \left( \sum_i A_i \right) B \end{aligned}$$

- So the NLL in terms of the precision matrix  $\Theta$  and sample covariance  $S$  is

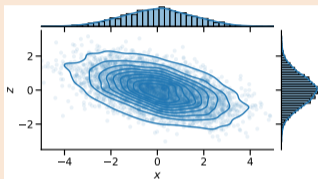
$$f(\Theta) = \frac{n}{2} \text{Tr}(S\Theta) - \frac{n}{2} \log |\Theta|, \text{ with } S = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \hat{\mu})(x^{(i)} - \hat{\mu})^T$$

- Weird-looking but has nice properties:
  - $\text{Tr}(S\Theta)$  is linear function of  $\Theta$ , with  $\nabla_{\Theta} \text{Tr}(S\Theta) = S$   
(it's the matrix version of an inner product  $s^T \theta$ ; called "Frobenius inner product")
  - Negative log-determinant is strictly convex, and  $\nabla_{\Theta} \log |\Theta| = \Theta^{-1}$   
(generalizes  $\nabla \log |x| = 1/x$  for  $x > 0$ )
- Using these two properties the **gradient matrix** has a simple form:

$$\nabla f(\Theta) = \frac{n}{2}(S - \Theta^{-1})$$

which is what **we use to get the MLE**

- With continuous features, we could model  $p(x, y)$  as a multivariate Gaussian



- Training could use the closed-form MLE/MAP for multivariate Gaussian
- We obtain a univariate Gaussian  $p(y | x)$  using conditioning formula,

$$y | x \sim \mathcal{N} \left( \mu_y + \Sigma_{yx} \Sigma_x^{-1} (x - \mu_x), \sigma_y^2 - \Sigma_{yx} \Sigma_x^{-1} \Sigma_{yx}^T \right)$$

- The conditional mean is a linear function,  $w^T x + b$
- Could extend to multiple outputs, with correlations given based on  $\Sigma_y$