CPSC 440: Advanced Machine Learning Learning Gaussians

Mark Schmidt

University of British Columbia

Winter 2022

Last Time: Inference in Multivariate Gaussian

 ${\ensuremath{\bullet}}$ The multivariate normal/Gaussian distribution models PDF of vector x^i as

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

where $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathbb{R}^{d \times d}$ is symmetric with $\Sigma \succ 0.$

- The density for a linear transformation of a product of independent Gaussians.
- Models correlations between non-zero elements of Σ .
 - If $\boldsymbol{\Sigma}$ is diagonal then it assumes all variables are independent.
- We discussed affine property that linear transformation of Gaussians is Gaussian.
 - Can be used to generate samples from a Gaussian.
- We discussed how marginals and conditionals are Gaussian.
- We often draw a graph of the non-zero elements of precision $\Theta = \Sigma^{-1]}$.
 - Variables are conditionally independent if conditioning set blocks paths in the graph.

Discussion of Independence in Gaussians

- If Σ is diagonal then Θ is diagonal.
 - This gives a disconnected graph: all variables are independent.
- If Θ is a full matrix, graph does not imply any conditional independences.
 - "Everything depends on everything, no matter how many of the x_j you know."
- The value Θ_{ij} is related to the partial correlation which is -Θ_{ij}/√Θ_{ii}Θ_{jj}.
 The "remaining correlation when we know all other variables".
- Dependencies can exist if $\Theta_{ij} = 0$ due to correlations with other variables.
 - Only independent if all paths that correlation could go across are blocked.

MLE for Multivariate Gaussian (Mean Vector)

• With a multivariate Gaussian we have

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

so up to a constant our negative log-likelihood for n examples x^i is

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

• This is a convex quadratic in μ , setting gradient to zero gives

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

• MLE for μ is the mean along each dimension, and it does not depend on Σ .

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

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

• After some tedious linear algebra (in bonus slides) we obtain that this is equal to

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

where:

- S is the covariance of the data (if we subtract mean from all examples, $S = (1/n)X^TX$).
- Trace operator Tr(A) is the sum of the diagonal elements of A.

• Gradient matrix of NLL with respect to Θ is (not obviously)

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

• The MLE for a given μ is obtained by setting gradient matrix to zero, giving

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

- The constraint $\Sigma \succ 0$ means we need positive-definite sample covariance, $S \succ 0$.
 - $\bullet~$ If S is not positive-definite, NLL is unbounded below and no MLE exists.
 - This is like requiring "not all values are the same" in univariate Gaussian.
 - In d-dimensions, you need d linearly-independent x^i values (no "collinearity")
- For most distributions, the MLEs are not the data's mean and covariance.

MAP Estimation for Mean

• For fixed Σ , conjugate prior for mean is a Gaussian:

$$x^i \sim \mathcal{N}(\mu, \Sigma), \mu \sim \mathcal{N}(\mu_0, \Sigma_0), \rightarrow \mu \mid X, \Sigma \sim \mathcal{N}(\mu^+, \Sigma^+),$$

where (using product of Gaussians property we are about to cover)

$$\begin{split} \Sigma^+ &= (n\Sigma^{-1} + \Sigma_0^{-1})^{-1}, \\ \mu^+ &= \Sigma^+ (n\Sigma^{-1}\mu_{\mathsf{MLE}} + \Sigma_0^{-1}\mu_0). \end{split} \qquad \text{MAP estimate of } \mu \end{split}$$

• In special case of $\Sigma=\sigma^2 I$ and $\Sigma_0=(1/\lambda)I$ we get

$$\Sigma^{+} = ((n/\sigma^{2})I + \lambda I)^{-1},$$

$$\mu^{+} = \Sigma^{+}((n/\sigma^{2})\mu_{\mathsf{MLE}} + \lambda\mu_{0}).$$

- Posterior predictive is $\mathcal{N}(\mu^+, \Sigma + \Sigma^+)$ (product of (n+2) then marginalize).
 - Many Bayesian inference tasks have closed form, or Monte Carlo is easy.

Product of Gaussian Densities Property

• Consider variable x whose PDF is written as product of two Gaussians,

 $p(x) = f_1(x)f_2(x)$

where:

- f_1 is proportional to a Gaussian with mean μ_1 and covariance I.
- f_2 is proportional to a Gaussian with mean μ_2 and covariance I.



Product of Gaussian Densities Property

- If $p(x) \propto f_1(x) f_2(x)$ f_1 and f_2 with
 - f_1 proportional to a Gaussian with mean μ_1 and covariance Σ_1 .
 - f_2 proportional to a Gaussian with mean μ_2 and covariance Σ_2 .
- Then p is a Gaussian with (see textbook)

covariance of
$$\Sigma = (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}$$
.
mean of $\mu = \Sigma \Sigma_1^{-1} \mu_1 + \Sigma \Sigma_2^{-1} \mu_2$,

• How we do we use this to derive the posterior distribution for the mean?

$$p(\mu \mid X, \Sigma, \mu_0, \Sigma_0) \propto p(\mu \mid \mu_0, \Sigma_0) \prod_{i=1}^n p(x^i \mid \mu, \Sigma)$$
(Bayes rule)
$$= p(\mu \mid \mu_0, \Sigma_0) \prod_{i=1}^n p(\mu \mid x^i, \Sigma)$$
(symmetry of x^i and μ)
$$= (\text{product of } (n+1) \text{ Gaussians}).$$

MAP Estimation in Multivariate Gaussian (Trace Regularization)

 \bullet A common MAP estimate for Σ is

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

where S is the covariance of the data.

• Key advantage: $\hat{\Sigma}$ is postiive-definite (eigenvalues are at least λ).

• This minimizes NLL plus L1-regularization of precision diagonals (see bonus)

$$f(\Theta) = \underbrace{\mathsf{Tr}(S\Theta) - \log |\Theta|}_{\mathsf{NLL}} + \lambda \sum_{j=1}^{d} |\Theta_{jj}|.$$

although it does not set Θ_{ii} values to exactly zero.

• Log-determinant term becomes arbitrarily steep as the Θ_{jj} approach 0.

Graphical LASSO

• A popular generalization called the graphical LASSO,

$$f(\Theta) = \mathsf{Tr}(S\Theta) - \log|\Theta| + \lambda \sum_{i=1}^{d} \sum_{j=1}^{d} |\Theta_{ij}|,$$

where we apply L1-regularization to all elements of Θ .

- With large enough λ , gives sparse off-diagonals in Θ .
 - Though need specialized optimization algorithms to solve this problem.
- Recall that sparsity of Θ determines conditional independence.
 - When we set a $\Theta_{ij} = 0$ it remove an edges from the graph.
 - Makes the graph simpler, and can make computations cheaper.

Learning in Multivariate Gaussians

Supervised Learning with Gaussians

Graphical LASSO Example

• Graphical LASSO applied to stocks data:



https://normaldeviate.wordpress.com/2012/09/17/high-dimensional-undirected-graphical-models

Graphical LASSO Example

• Graphical LASSO applied to US senate voting data (Bush junior era):



 ${\tt https://normaldeviate.wordpress.com/2012/09/17/high-dimensional-undirected-graphical-models}$

Graphical LASSO Example

• Graphical LASSO applied to protein data:



https://normaldeviate.wordpress.com/2012/09/17/high-dimensional-undirected-graphical-models

Graphical LASSO on Digits

• Precision matrix from graphical LASSO applied to MNIST digits ($\lambda = 1/8$):



- To understand this picture, first the size of the precision matrix:
 - The images of digits, which are m imes m matrices (m pixels by m pixels)
 - This gives $d = m^2$ elements of x^i , which we'll assume are in "column-major" order.
 - $\bullet\,$ Frist m elements of x^i are column 1, next m elements are columm 2, and so on.
 - The picture above, which is $d \times d$ so will thus be $m^2 \times m^2$.

Graphical LASSO on Digits

• Precision matrix from graphical LASSO applied to MNIST digits ($\lambda = 1/8$):



- So what are the non-zeroes in the precision matrix?
 - **(**) The diagonals $\Theta_{i,i}$ (positive-definite matrices must have positive diagonals).
 - 2 The first off-diagonals $\Theta_{i,i+1}$ and $\Theta_{i+1,i}$.
 - This represents the dependencies between adjacent pixels vertically.
 - **③** The (m+1) off-diagonals $\Theta_{i,i+m}$ and $\Theta_{i+m,i}$.
 - This represents the dependencies between adjacent pixels horizontally.
 - Because in "column-major" order, you go "right" a pixel every m indices.

Graphical LASSO on Digits

• Precision matrix from graphical LASSO applied to MNIST digits ($\lambda = 1/8$):



- The edges in the graph are pixels next to each other in the image.
- Graphical Lasso is a special case of structure learning in graphical models.
 We will discuss graphical models more later.

Conjugate Priors for Covariance

- Graphical LASSO is not using a conjugate prior.
- \bullet Conjugate prior for Θ 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 Σ is inverse-Wishart (equivalent posterior).
- If both μ and Θ 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

Learning in Multivariate Gaussians

Supervised Learning with Gaussians

Outline



2 Supervised Learning with Gaussians

Generative Classification with Gaussians

• We previously considerd the generative classifier naive Bayes.

- Assumed $x_i \perp x_j \mid y$, which is strong/unrealistic.
- Consider a generative classifier with continuous features:

$$p(y^{i} \mid x^{i}) \propto p(x^{i}, y^{i})$$

$$= \underbrace{p(x^{i} \mid y^{i})}_{\text{continuous discrete}} \underbrace{p(y^{i})}_{\text{discrete}}.$$

- In Gaussian discriminant analysis (GDA) we assume $x^i \mid y^i$ is a Gaussian.
 - It is classification so output y^i is categorical.
 - Classifier asks "which Gaussian makes x^i most likely?"
 - This can model pairwise correlations within each class.
 - Does not need naive Bayes assumption.

Gaussian Discriminant Analysis (GDA) and Closed-Form MLE

• In Gaussian discriminant analysis we assume $x^i \mid y^i$ is a Gaussian.

$$p(x^i, y^i = c) = \underbrace{p(y^i)p(x^i \mid y^i = c)}_{\text{product rule}} = \underbrace{\pi_c}_{p(y^i = c)} \underbrace{p(x^i \mid \mu_c, \Sigma_c)}_{\text{Gaussian PDF}}.$$

- A special case is linear discriminant analysis (LDA):
 - In LDA we assume that Σ_c is the same for all classes c.
- In LDA the MLE has a simple closed-form expression:

$$\hat{\pi}_c = \frac{n_c}{n}, \quad \hat{\mu}_c = \frac{1}{n_c} \sum_{y^i = c} x^i.$$

• $\hat{\pi}_c$ is fraction of times we are in class c, $\hat{\mu}$ is mean of class c.

Linear Discriminant Analysis (LDA)

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



https://web.stanford.edu/~hastie/Papers/ESLII.pdf

- LDA is a linear classifier.
 - Unlike other linear classifiers like logistic regression, it has a closed-form MLE.
 - Though it may be less accurate if the classes do not look like Gaussians.
- If class proportions π_c are equal, class label is determined by nearest mean.
 - Prediction is like a "1-nearest neighbour" or k-means clustering method.

Gaussian Discriminant Analysis (GDA)

- $\bullet\,$ We can also have a covariance Σ_c for each class.
 - So the class will be determined by class proportions, means, and variances.
- The MLE for each each Σ_c is the covariance of data in class c,

$$\hat{\Sigma}_c = \frac{1}{n_c} \sum_{y^i = c} (x_i - \hat{\mu}_c) (x_i - \hat{\mu}_c)^T,$$



https://web.stanford.edu/~hastie/Papers/ESLII.pdf

- This leads to a quadratic classifier.
 - GDA is sometimes called quadratic discriminant analysis.

Regression with Gaussians

• Regression is a variant on supervised learning where y^i is continuous.



https://en.wikipedia.org/wiki/Regression_analysis

- It is possible to use generative regression models.
 - For example, we could model p(x, y) as a multivariate Gaussian.
 - Then use that the conditional $p(y \mid x)$ is Gaussian for prediction.
- But we usually treat features as fixed (as in discriminative classification models).
 And to start, we will consider models that make linear predictions, ŷⁱ = w^Txⁱ.

Review: L2-Regularized Least Squares and Gaussians

• A common linear regression model is L2-regularized least squares,

$$\underset{w}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|Xw-y\|^2 + \frac{\lambda}{2} \|w\|^2,$$

• This corresponds to MAP estimation with a Gaussian likelihood and prior,

$$y^i \sim \mathcal{N}(w^T x^i, \sigma^2), \quad w_j \sim \mathcal{N}(0, \lambda^{-1}).$$

• By setting the gradient to zero, the unique solution is given by:

$$\hat{w} = \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} X^T X + \lambda I \right)^{-1} X^T y.$$

- In 340 we fixed $\sigma^2 = 1$ (since changing is σ^2 equivalent to changing λ). • In Bayesian inference, both σ^2 and λ affect the predictions.
- To predict on new example \tilde{x} with MAP estimate, we use $\hat{y} = \hat{w}^T \tilde{x}$.

Summary

- MLE for multivariate Gaussian:
 - MLE for μ is mean of data, MLE for Σ is covariance of data (if positive definite).
- Posterior and posterior predictive under Gaussian prior on mean is Gaussian.
 - Can be shown using that product of Gaussians is Gaussian.
- Graphical Lasso uses L1-regularization of precision matrix.
 - Leads to a sparse graph structure representing conditional independences.
- Supervised learning with Gaussians
 - Generative classifier with Gaussian classes is Gaussian discriminant analysis (GDA).
 - L2-regularized least squares is obtained using a Gaussian likelihood and prior.
 - Regression model assuming features fixed/non-random as in discriminative classifiers.
- Next time: linear regression plus empirical Bayes to cheat on your stats homework.

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

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

• Where the trace Tr(A) is the sum of the diagonal elements of A.

• That Tr(ABC) = Tr(CAB) when dimensions match is the cyclic property of trace.

 \bullet From the last slide we have in terms of precision matrix Θ that

$$= \frac{1}{2} \sum_{i=1}^{n} \operatorname{Tr}((x^{i} - \mu)(x^{i} - \mu)^{\top} \Theta) - \frac{n}{2} \log |\Theta|$$

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

$$=\frac{1}{2}\operatorname{Tr}\left(\sum_{i=1}^{n} (x^{i} - \mu)(x^{i} - \mu)^{\top}\Theta\right) - \frac{n}{2}\log|\Theta| \qquad \sum_{i}\operatorname{Tr}(A_{i}B) = \operatorname{Tr}\left(\sum_{i}A_{i}B\right)$$
$$=\frac{n}{2}\operatorname{Tr}\left(\left(\underbrace{\frac{1}{n}\sum_{i=1}^{n} (x^{i} - \mu)(x^{i} - \mu)^{\top}}_{\text{sample covariance 'S'}}\right)\Theta\right) - \frac{n}{2}\log|\Theta|. \qquad \left(\sum_{i}A_{i}B\right) = \left(\sum_{i}A_{i}\right)B$$

 $\bullet\,$ So the NLL in terms of the precision matrix Θ and sample covariance S is

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

- Weird-looking but has nice properties:
 - $\operatorname{Tr}(S\Theta)$ is linear function of Θ , with $\nabla_{\Theta} \operatorname{Tr}(S\Theta) = S$.

(it's the matrix version of an inner-product $s^{\top}\theta$) • Negative log-determinant is strictly-convex and has $\nabla_{\Theta} \log |\Theta| = \Theta^{-1}$.

(generalizes $\nabla \log |x| = 1/x$ for for x > 0).

• Using these two properties the gradient matrix has a simple form:

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

Trace Regularization and L1-regularization

 \bullet A classic regularizer for Σ is to add a diagonal matrix to S and use

 $\Sigma = S + \lambda I,$

which satisfies $\Sigma \succ 0$ because $S \succeq 0$ (eigenvalues at least λ).

• This corresponds to L1-regularization of diagonals of precision.

$$\begin{split} f(\Theta) &= \operatorname{Tr}(S\Theta) - \log |\Theta| + \lambda \sum_{j=1}^{d} |\Theta_{jj}| & (\text{Gauss. NLL plus L1 of diags}) \\ &= \operatorname{Tr}(S\Theta) - \log |\Theta| + \lambda \sum_{j=1}^{d} \Theta_{jj} & (\text{Diagonals of pos. def. matrix are } > 0) \\ &= \operatorname{Tr}(S\Theta) - \log |\Theta| + \lambda \operatorname{Tr}(\Theta) & (\text{Definition of trace}) \\ &= \operatorname{Tr}(S\Theta + \lambda\Theta) - \log |\Theta| & (\text{Linearity of trace}) \\ &= \operatorname{Tr}((S + \lambda I)\Theta) - \log |\Theta| & (\text{Distributive law}) \end{split}$$

- Taking gradient and setting to zero gives $\Sigma = S + \lambda$.
 - But doesn't set to exactly zero as log-determinant term is too "steep" at 0.