

Solution 1

a)

$$X = (A^T Q^{-1} A)^{-1} A^T Q^{-1} Y$$

b)

$$L(X, \lambda) = \frac{1}{2} (Y - AX)^T Q^{-1} (Y - AX) + \lambda (b^T X - c)$$

$$\frac{\partial}{\partial X} L(X, \lambda) = -(Y - AX)^T Q^{-1} A + \lambda b^T = 0$$

$$\frac{\partial}{\partial \lambda} L(X, \lambda) = b^T X - c = 0$$

$$X = (A^T Q^{-1} A)^{-1} (A^T Q^{-1} Y - \frac{1}{2} \lambda b)$$

where λ is the Lagrange multiplier. Then by solving the two equations to get X .

c)

$$\mathcal{L}(\mathbf{X}, \lambda_1, \lambda_2) = (\mathbf{Y} - \mathbf{A}\mathbf{X})^T \mathbf{Q}^{-1} (\mathbf{Y} - \mathbf{A}\mathbf{X}) + \lambda_1 (\mathbf{b}^T \mathbf{X} - c) + \lambda_2 (\mathbf{X}^T \mathbf{X} - d)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}} = -2\mathbf{A}^T \mathbf{Q}^{-1} (\mathbf{Y} - \mathbf{A}\mathbf{X}) + \lambda_1 \mathbf{b} + 2\lambda_2 \mathbf{X} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_1} = \mathbf{b}^T \mathbf{X} - c = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_2} = \mathbf{X}^T \mathbf{X} - d = 0$$

Then by solving the two equations to get X .

d)

If both A and X are unknown, solving them requires an iterative method due to the joint estimation problem. This can be approached by using alternating minimization.

First, to solve for X with a fixed A ,

$$\mathcal{L}(X, \lambda) = (Y - AX)^T Q^{-1} (Y - AX) + \lambda (X^T X - d)$$

$$\frac{\partial \mathcal{L}}{\partial X} = -2A^T Q^{-1} (Y - AX) + 2\lambda X = 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = X^T X - d = 0$$

Solving these equations will yield the value for X .

Next, fix X and solve for A ,

$$\begin{aligned}\mathcal{L}(A, \lambda') &= (Y - AX)^\top Q^{-1}(Y - AX) + \lambda'(\text{Trace}(A^\top A) - e) \\ \frac{\partial \mathcal{L}}{\partial A} &= -2Q^{-1}(Y - AX)X^\top + 2\lambda' A = 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda'} &= \text{Trace}(A^\top A) - e = 0\end{aligned}$$

Solving these equations will yield the value for A .

By iterating between these two steps, adjusting X with a fixed A and then adjusting A with the new X , and repeating until the change in the values of X and A falls below a certain threshold, we find the optimal solution to the least squares problem with the given constraints.

Solution 2

Conditional distribution

$$p(Y|X) = \mathcal{N}(\mathbf{Y}|AX, \beta^{-1}\mathbf{I})$$

Joint distribution

$$p(Y, X) = \mathcal{N}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \middle| \begin{bmatrix} m_0 \\ Am_0 \end{bmatrix}, \begin{bmatrix} \Sigma_0 & \Sigma_0 A^\top \\ A\Sigma_0 & A\Sigma_0 A^\top + \beta^{-1}I \end{bmatrix}\right)$$

Marginal distribution

$$p(Y) = \int p(Y, X) dX = \mathcal{N}(Y|Am_0, A\Sigma_0 A^\top + \beta^{-1}I)$$

Posterior distribution

$$\begin{aligned}p(X|Y = \mathbf{y}, \beta, \mathbf{m}_0, \Sigma_0) &\propto p(Y = \mathbf{y}|X, \beta)p(X|\mathbf{m}_0, \Sigma_0) \\ &= \mathcal{N}(X|\mu_{X|Y}, \Sigma_{X|Y})\end{aligned}$$

$$\begin{aligned}\mu_{X|Y} &= \mu_X + \Sigma_{XY}\Sigma_{YY}^{-1}(\mathbf{y} - \mu_Y) = m_0 + \Sigma_0 A^\top (A\Sigma_0 A^\top + \beta^{-1}I)^{-1}(\mathbf{y} - Am_0) \\ \Sigma_{X|Y} &= \Sigma_{XX} - \Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX} = \Sigma_0 - \Sigma_0 A^\top (A\Sigma_0 A^\top + \beta^{-1}I)^{-1}A\Sigma_0\end{aligned}$$

Posterior predictive distribution

$$\begin{aligned}p(\tilde{Y}|Y = \mathbf{y}, \beta, m_0, \Sigma_0) &= \int p(\tilde{Y}|X)p(X|Y = \mathbf{y}, \beta, m_0, \Sigma_0) dX \\ &= \mathcal{N}(X|\mu_{X|Y}, \Sigma_{X|Y}) \times \mathcal{N}(\mathbf{Y}|AX, \beta^{-1}\mathbf{I})\end{aligned}$$

Prior predictive distribution

$$p(Y|\beta, m_0, \Sigma_0) = \mathcal{N}(Y|Am_0, A\Sigma_0 A^\top + \beta^{-1}I)$$

Solution 3

The posterior distribution is

$$p(\mathbf{w}|\mathcal{D}, \beta, \mathbf{m}_0, \alpha) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

where

$$\begin{aligned}\mathbf{m}_N &= \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\Phi^T\mathbf{t}) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta\Phi^T\Phi \\ \mathbf{S}_0 &= \alpha^{-1}\mathbf{I}\end{aligned}$$

The posterior predictive distribution is

$$p(\hat{y}|\hat{x}, \mathcal{D}, \beta, \mathbf{m}_0, \alpha) = \mathcal{N}(\hat{y}|\mathbf{m}_N^T\phi(\hat{x}), \sigma_N^2(\hat{x}))$$

where

$$\sigma_N^2(\hat{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x})$$

The prior predictive distribution is

$$p(\mathcal{D}|\beta, \mathbf{m}_0, \alpha) = \prod_{n=1}^n \mathcal{N}(y_n|\phi_n^T \mathbf{m}_0, \phi_n^T \alpha^{-1} \phi_n + \beta^{-1})$$

Solution 4. Logistics Regression

The posterior distribution

$$p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \alpha) \propto p(\mathbf{w}|\mathbf{m}_0, \alpha)p(\mathcal{D}|\mathbf{w}) = p(\mathbf{w}|\mathbf{m}_0, \alpha) \prod_{n=1}^N p(t_n|\mathbf{w})$$

$$\ln p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \alpha) = -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{i=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + C$$

$$\mathbf{S}_N^{-1} = -\nabla \nabla \ln p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \alpha) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n)\phi_n\phi_n^T$$

Hence,

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{MAP}, \mathbf{S}_N)$$

The posterior predictive distribution

$$p(t|x, \mathcal{D}, \mathbf{m}_0, \alpha)$$

The prior predictive distribution

$$p(\mathcal{D}|\mathbf{m}_0, \alpha) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{m}_0, \alpha)$$

Solution 5

(1)

$$\frac{\partial y}{\partial w^{(1)}} = \frac{\partial y}{\partial a_2} \cdot \frac{\partial a_2}{\partial z} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial a_1}{\partial w^{(1)}} = y(1 - y)w^{(2)}h'(a_1)x$$

$$\begin{aligned}\frac{\partial y}{\partial w^{(2)}} &= \frac{\partial y}{\partial a_2} \cdot \frac{\partial a_2}{\partial w^{(2)}} = y(1-y)z \\ \frac{\partial y}{\partial a_1} &= \frac{\partial y}{\partial a_2} \cdot \frac{\partial a_2}{\partial z} \cdot \frac{\partial z}{\partial a_1} = y(1-y)w^{(2)}h'(a_1) \\ \frac{\partial y}{\partial a_2} &= \frac{\partial y}{\partial a_2} = y(1-y) \\ \frac{\partial y}{\partial x} &= \frac{\partial y}{\partial a_2} \cdot \frac{\partial a_2}{\partial z} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial a_1}{\partial x} = y(1-y)w^{(2)}h'(a_1)w^{(1)}\end{aligned}$$

(2)

Use MSE

$$\begin{aligned}L &= \frac{1}{2}(y - t)^2 \\ \Delta w^{(2)} &= -\eta \frac{\partial L}{\partial w^{(2)}} \\ &= -\eta \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_2} \frac{\partial a_2}{\partial w^{(2)}} \\ &= \eta(t - y)y'(a_2)z \\ &= \eta(t - y)y(1 - y)z \\ \Delta w^{(1)} &= -\eta \frac{\partial L}{\partial w^{(1)}} \\ &= -\eta \frac{\partial L}{\partial y} \frac{\partial y}{\partial a_2} \frac{\partial a_2}{\partial z} \frac{\partial z}{\partial a_1} \frac{\partial a_1}{\partial w^{(1)}} \\ &= \eta(t - y)y'(a_2)w^{(2)}h'(a_1)x \\ &= \eta(t - y)y(1 - y)w^{(2)}h'(a_1)x\end{aligned}$$

Solution 6

(a)

Posterior distribution

$$p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \alpha, \beta) \propto p(\mathbf{w}|\mathbf{m}_0, \alpha) \times p(D|\mathbf{w}, \beta)$$

where

$$\begin{aligned}p(\mathbf{w}|\mathbf{m}_0, \alpha) &= \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \alpha^{-1}I) \\ p(D|\mathbf{w}, \beta) &= \prod_{n=1}^N \mathcal{N}(t_n|y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})\end{aligned}$$

Posterior predictive distribution

$$p(t|x, D, \beta, \mathbf{m}_0, \alpha) = \int p(t|x, \mathbf{w})q(\mathbf{w}|D)d\mathbf{w} = \mathcal{N}(t|y(x, \mathbf{w}_{MAP}), \sigma^2(x))$$

where

$$\begin{aligned}y(x, \mathbf{w}) &\simeq y(x, \mathbf{w}_{MAP}) + g^T(\mathbf{w} - \mathbf{w}_{MAP}) \\ g &= \nabla_{\mathbf{w}} y(x, \mathbf{w})|_{\mathbf{w}=\mathbf{w}_{MAP}} \\ \sigma^2(x) &= \beta^{-1} + g^T A^{-1} g\end{aligned}$$

(b)

Posterior distribution

$$p(\mathbf{w}|D, \alpha) = \int p(t|x, \mathbf{w})q(\mathbf{w}|D, \alpha)d\mathbf{w}$$

Prior predictive distribution

$$p(D|\beta, \mathbf{m}_0, \alpha) = \prod_{n=1}^N p(t|x, \alpha)$$

Solution 7

a) Please explain why the dual problem formulation is used to solve the SVM machine learning problem.

The dual formulation of SVM is favored because it simplifies the problem when data isn't linearly separable. It allows the use of kernels to handle higher-dimensional space efficiently. The dual form guarantees finding a global minimum, as it's a convex problem. It's computationally efficient since only support vectors (a subset of data) determine the decision boundary, leading to a sparse and faster-to-compute model.

b)

- **i) SVM vs Logistic Regression:** SVM focuses on the widest margin between classes, using hinge loss which ignores errors outside the margin. Logistic regression considers all data points with a logistic loss, predicting probabilities.
- **ii) v-SVM vs Least Squares Regression:** v-SVM uses ϵ -insensitive loss, ignoring errors within a certain range, which makes it robust to outliers. Least squares regression minimizes the sum of squared errors, sensitive to all deviations, making it less robust to outliers.

c) Neural networks use logistic activation functions for several reasons:

Neural networks use logistic activation functions to introduce non-linearity, enabling the network to learn complex patterns. Logistic functions are useful in output layers for binary classification since they produce probabilities (values between 0 and 1). They're continuously differentiable, a property needed for training algorithms like backpropagation. However, due to issues like gradient saturation, other functions like ReLU are now preferred in hidden layers.

d)

- **Differences:**
 - The logistic (sigmoid) function outputs values between 0 and 1, which is ideal for binary classifications.
 - The ReLU (Rectified Linear Unit) function outputs zero for negative inputs and raw input for positive inputs, which helps with the vanishing gradient problem and speeds up training.
 - The tanh function outputs values between -1 and 1, which centers the data, improving the learning for subsequent layers.
- **Usage:**
 - Logistic functions are often used in the output layer for binary classification problems.
 - ReLU is preferred in hidden layers due to its efficiency and effectiveness in deep networks.

- Tanh is used when data centering is beneficial, but less common due to vanishing gradients.

e)

- The Jacobian matrix, representing first-order derivatives, is crucial for understanding the gradient of multivariate functions, helping in gradient descent optimization.
- The Hessian matrix, representing second-order derivatives, is used to find the curvature of the loss function, informing about the optimization landscape and guiding adjustments to the learning rate or direction.

f)

Exponential family distributions are common because they have convenient mathematical properties that allow for efficient estimation and inference, such as conjugate priors in Bayesian analysis.

Non-examples include the Cauchy distribution, which lacks a defined mean or variance, and the uniform distribution, which does not have a natural exponential form.

g)

Kullback-Leibler (KL) divergence measures how one probability distribution diverges from a second, expected probability distribution. It's useful for machine learning because it quantifies the difference between the learned model distribution and the true distribution of data.

- Example 1: In variational autoencoders (VAEs), KL divergence is used to regularize the encoder by penalizing divergences from the prior distribution, encouraging the latent space to approximate a standard normal distribution.
- Example 2: In natural language processing, KL divergence helps in comparing the similarity of word distribution in different text documents, aiding in tasks like topic modeling.

h)

Data augmentation techniques are a form of regularization for neural networks because they artificially increase the diversity of data available for training. By applying transformations like rotation, scaling, or cropping, they help the model generalize better, reduce overfitting, and improve robustness to variations in new, unseen data.

i)

- **Central Limit Theorem:** Many natural phenomena tend to follow a Gaussian distribution when they are the sum of many independent random variables.
- **Conjugacy:** Gaussian distributions are mathematically tractable, especially as conjugate priors in Bayesian inference, simplifying the computation of posterior distributions.
- **Continuity and Differentiability:** Gaussians are smooth and differentiable, which is beneficial for optimization algorithms that rely on gradient information.
- **Descriptive:** Gaussian distributions are defined by just two parameters (mean and variance), which can effectively capture the characteristics of many real-world datasets.

j)

it simplifies the computation by approximating the distribution around the peak (mode) with a Gaussian. This is effective when the peak is sharp, as it captures the main contribution to the integral, making it a practical approach for high-dimensional problems where exact computation is infeasible.

k)

Balance fit and complexity to prevent underfitting and overfitting. Cross-validation, information criteria (AIC, BIC), and regularization are used to evaluate and select models.

l)

features should be relevant, non-redundant, and have predictive power. Techniques like feature importance, correlation analysis, and domain knowledge can guide selection. For testing, samples should be representative but unseen during training to provide an unbiased evaluation.

An example is splitting a dataset into training and test sets, ensuring the test set includes a variety of examples across the feature space.

m)

The MAP model is often preferred over the ML model because it incorporates prior knowledge about the parameters through the prior distribution. This can lead to more reliable estimates, especially with small datasets. MAP can also mitigate overfitting by penalizing complex models, whereas ML estimates can overfit by focusing solely on the data likelihood.

Solution 8

(1)

Generative approaches to machine learning model how the data is generated, by learning the joint probability distribution $P(x, y)$. They can generate new data points and are powerful in unsupervised learning tasks. However, they can be complex and computationally intensive because they learn the full data distribution.

Discriminative approaches model the decision boundary between classes directly by learning the conditional probability $P(y|x)$. They often require less computation and provide better performance on classification tasks.

Example: Consider spam email filtering. A generative model would learn the distribution of words in both spam and non-spam emails to classify or even generate new email content. A discriminative model would learn the boundary that separates spam from non-spam based on features of the emails.

(2) Analyzing the GAN Model:

GAN combine both generative and discriminative models.

- The generative model in a GAN learns to produce data that's indistinguishable from real data, aiming to capture the data distribution.
- The discriminative model learns to distinguish real data from the fake data produced by the generator.

