

Assignment Project Exam Help

Machine learning lecture slides

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Multivariate Gaussians and PCA

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- ▶ Multivariate Gaussians
- ▶ Eigendecompositions and covariance matrices
- ▶ Principal component analysis
- ▶ Principal component regression and spectral regularization
- ▶
- ▶

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Multivariate Gaussians: Isotropic Gaussians

- ▶ Start with $X = (X_1, \dots, X_d) \sim N(0, I)$, i.e., X_1, \dots, X_d are iid $N(0, 1)$ random variables.

Assignment Project Exam Help

- ▶ Probability density function is product of (univariate) Gaussian densities

- ▶ $(X_i) = 0$

$/ j$
<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Assignment Project Exam Help
<https://eduassistpro.github.io>
Add WeChat edu_assist_pr

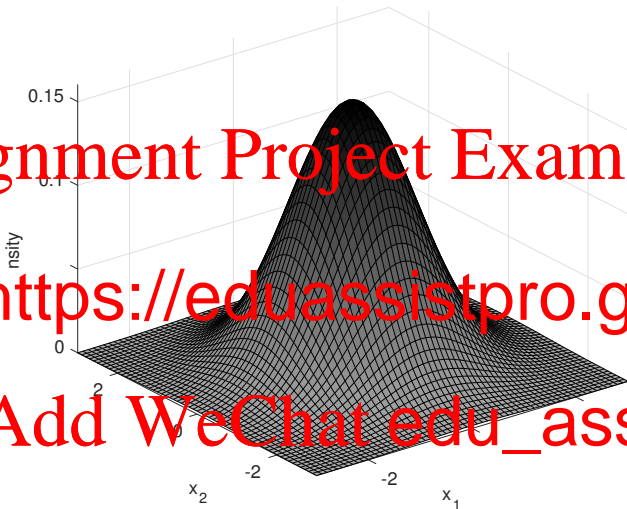


Figure 1: Density function for isotropic Gaussian in \mathbb{R}^2

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

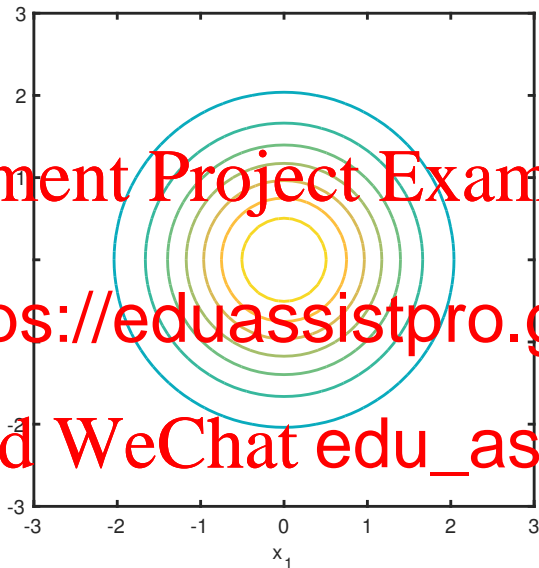


Figure 2: Density function level sets for isotropic Gaussian in \mathbb{R}^2

Affine transformations of random vectors

- ▶ Start with any random vector Z , then apply linear transformation, followed by translation

▶ $X := MZ + \mu$, for $M \in \mathbb{R}^k \times \mathbb{R}^d$ and $\mu \in \mathbb{R}^k$

▶ Fact: $\mathbb{E}(X) = M\mathbb{E}(Z) + \mu$, $\text{cov}(X) = M\text{cov}(Z)M^T$

- ▶ E.g., let $u \in \mathbb{R}^d$ be a unit vector ($u^T u = 1$), and $X := u^T Z$ (a scalar), and

- ▶ <https://eduassistpro.github.io>

distribution, not just Gaussian distribution

- ▶ However, it is convenient to illustrate the transformations of Gaussian distribution

the Gaussian pdf is easy to understand.

Multivariate Gaussians: General Gaussians

- ▶ If $Z \sim \mathcal{N}(0, I)$ and $X = MZ + \mu$, we have $\mathbb{E}(X) = \mu$ and $\text{cov}(X) = MM^T$

Assignment Project Exam Help

- ▶ Assume $M \in \mathbb{R}^{d \times d}$ is invertible (else we get a degenerate Gaussian distribution).
- ▶ We say $X \sim \mathcal{N}(\mu, MM^T)$

<https://eduassistpro.github.io>

- ▶ Note: every non-singular covariance mat MM^T for some non-singular matrix

Add WeChat edu_assist_pro

Assignment Project Exam Help
<https://eduassistpro.github.io>
Add WeChat edu_assist_pr

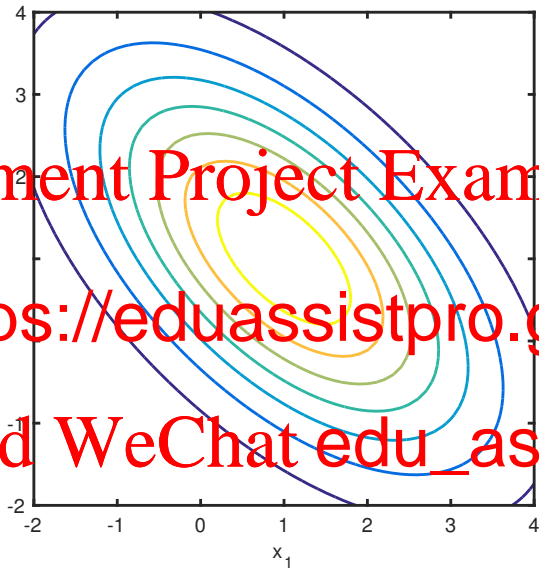


Figure 3: Density function level sets for anisotropic Gaussian in \mathbb{R}^2

Inference with multivariate Gaussians (2)

- ▶ Bivariate case: $(X_1, X_2) \sim \mathcal{N}(\mu, \Sigma)$ in \mathbb{R}^2

Assignment Project Exam Help

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{1,1} & \Sigma_{1,2} \\ \Sigma_{2,1} & \Sigma_{2,2} \end{bmatrix}$$


<https://eduassistpro.github.io>

- ▶ Miracle 1: it is a Gaussian distribution
- ▶ Miracle 2: mean provided by linear prediction with smallest MSE
- ▶ Miracle 3: variance doesn't depend on

Add WeChat edu_assist_pro

Inference with multivariate Gaussians (2)

- ▶ What is the distribution of $X_2 \mid X_1 = x_1$?
 - ▶ Miracle 1: it is a Gaussian distribution
 - ▶ Miracle 2: mean provided by linear prediction of X_2 from X_1 with smallest MSE
 - ▶ Miracle 3: variance doesn't depend on x_1

e:

<https://eduassistpro.github.io>

$$\hat{m} = \frac{X_1, X_2}{\text{var}(X_1)} \frac{\Sigma_{1,2}}$$

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

- ▶ Therefore:

$$\begin{aligned}\mathbb{E}[X_2 \mid X_1 = x_1] &= \hat{m}x_1 + \hat{\theta} \\ &= \mu_2 + \hat{m}(x_1 - \mu_1) \\ &= \mu_2 + \frac{\Sigma_{1,2}}{\Sigma_{1,1}}(x_1 - \mu_1)\end{aligned}$$

Inference with multivariate Gaussians (3)

- ▶ What is the distribution of $X_2 \mid X_1 = x_1$?

- ▶ Miracle 1: it is a Gaussian distribution

- ▶ Miracle 2: mean provided by linear prediction of X_2 from X_1 with smallest MSE

- ▶ Miracle 3: variance doesn't depend on x_1

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

$$= \Sigma_{2,2}$$

$$= \Sigma_{2,2}$$

$$= \Sigma_{2,2} - \frac{\Sigma_{1,2}^2}{\Sigma_{1,1}}$$

$$= \Sigma_{2,2} - \frac{\Sigma_{1,2}^2}{\Sigma_{1,1}}$$

- Beyond bivariate Gaussians: same as above, but just writing things properly using matrix notations

Assignment Project Exam Help

$$\mathbb{E}[X_2 | X_1 = x_1] = \mu_2 + \Sigma_{2,1} \Sigma_{1,1}^{-1} (x_1 - \mu_1)$$

1

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Eigendecomposition (1)

- ▶ Every symmetric matrix $M \in \mathbb{R}^{d \times d}$ has d real eigenvalues, which we arrange as

Assignment Project Exam Help
 $\lambda_1 \geq \dots \geq \lambda_d$



<https://eduassistpro.github.io>

- ▶ This means

Add WeChat $M v_i$ edu_assist_pr
for each $i = 1, \dots, d$, and

$$v_i^\top v_j = \mathbf{1}_{\{i=j\}}$$

Eigendecomposition (2)

- ▶ Arrange v_1, \dots, v_d in an orthogonal matrix $V := [v_1 | \dots | v_d]$

- ▶ $V^T V = I$ and $V V^T = \sum_{i=1}^d v_i v_i^T = I$

- ▶ Therefore,

$$M = M V V^T$$

d

Assignment Project Exam Help

<https://eduassistpro.github.io>

$$\lambda_i v_i v_i^T$$

$i=$

- ▶ This is our preferred way to express the orthogonal matrix

- ▶ Also called spectral decomposition

- ▶ Can also write $M = V \Lambda V^T$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$

- ▶ The matrix V diagonalizes M :

$$V^T M V = \Lambda$$

Covariance matrix (1)

► $A \in \mathbb{R}^{n \times d}$ is data matrix

► $\Sigma := A^T A = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ is

(empirical) second-moment matrix

► If $\frac{1}{n} \sum_{i=1}^n x_i = 0$ (data are “centered”), this is the

► <https://eduassistpro.github.io>

Add WeChat $\frac{1}{n} \sum_{i=1}^n x_i x_i^T$ edu_assist_pro

is (empirical) variance of data along direction u

Covariance matrix (2)

- ▶ Note: some pixels in OCR data have very little (or zero!) variation

- ▶ These are “coordinate directions” (e.g. $u = (1, 0, \dots, 0)$)
- ▶ Probably can/should ignore these!

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Figure 4: Which pixels are likely to have very little variance?

Top eigenvector

- ▶ Σ is symmetric, so can write eigendecomposition

Assignment Project Exam Help

$$\Sigma = \sum_{i=1}^d \lambda_i v_i v_i^\top$$

- ▶
- ▶ <https://eduassistpro.github.io>

- ▶ This follows from the following charact

Add WeChat edu_assist_pr

$$\frac{1}{\sqrt{\lambda_1}} \Sigma v_1 = \max_{u \in \mathbb{R}^d, \|u\|=1}$$

Assignment Project Exam Help
<https://eduassistpro.github.io>
Add WeChat edu_assist_pr

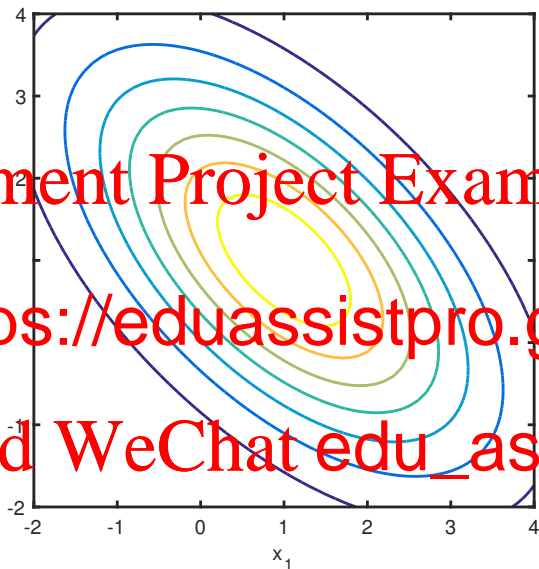


Figure 5: What is the direction of the top eigenvector for the covariance of this Gaussian?

Top k eigenvectors

- ▶ What about among directions orthogonal to v_1 ?
 - ▶ Answer: v_2 , corresponding to second largest eigenvalue λ_2
- ▶ (Note: all eigenvalues of L are non-negative!)
- ▶ For any k , $V_k := [v_1 | \cdots | v_k]$ satisfies

<https://eduassistpro.github.io> $\sum_{i=1}^k \lambda_i$

(the top k eigenvectors)

Add WeChat edu_assist_pr

Principal component analysis

- ▶ k -dimensional principal components analysis (PCA) mapping:

Assignment Project Exam Help

$$\varphi(x) = (x \cdot v_1, \dots, x \cdot v_k)^T = V_k^T x, x \in \mathbb{R}^d$$

where $V = [v_1 \quad \dots \quad v_k] \in \mathbb{R}^{d \times k}$

- ▶
- ▶ <https://eduassistpro.github.io>

Add WeChat edu_assist_pr

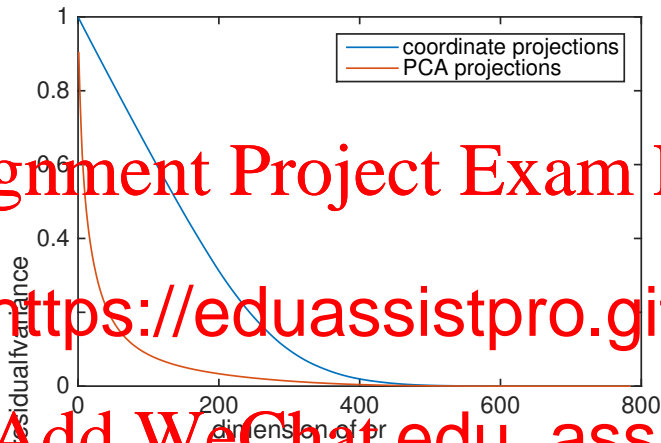


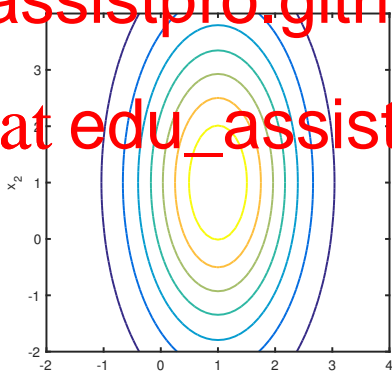
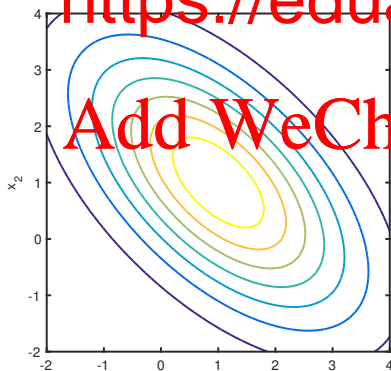
Figure 6: Fraction of residual variance from projection dimension

Covariance of data upon PCA mapping

- Covariance of data upon PCA mapping:

$$\frac{1}{n} \sum_{i=1}^n \langle \varphi(x_i), \varphi(x_i) \rangle = \frac{1}{n} \sum_{i=1}^n V_k^T x_i x_i^T V_k = V_k^T \Sigma V_k = \Lambda_k$$

where Λ_k is diagonal matrix with $\lambda_1, \dots, \lambda_k$ along diagonal.



Principal component regression

- ▶ Use $\hat{\beta} = \Lambda_k^{-1} V_k^T A^T b$ to predict on new $x \in \mathbb{R}^d$:

$$\begin{aligned} \hat{\beta}(x) &= (V_k^T x)^T \Lambda_k^{-1} V_k^T A^T b \\ &= x^T (V_k \Lambda_k^{-1} V_k^T) (A^T b) \end{aligned}$$

- ▶ <https://eduassistpro.github.io>

$$\hat{w} := (V_k \Lambda_k^{-1} V_k^T)^T$$

- ▶ This is called principal component regression (k is hyperparameter)
- ▶ Alternative hyper-parameterization: $\lambda > 0$; same as before but using the largest k such that $\lambda_k \geq \lambda$.

Spectral regularization

- ▶ PCR and ridge regression are examples of *spectral regularization*.

- ▶ For a function $g: \mathbb{R} \rightarrow \mathbb{R}$, write $g(M)$ to mean

<https://eduassistpro.github.io>

- ▶ I.e., g is applied to eigenvalues of
- ▶ Generalizes effect of polynomials: e.g.

$$M^2 = (V \Lambda V^T)(V \Lambda V^T)$$

- ▶ **Claim:** Can write each of PCR and ridge regression as

$$\hat{w} = g(A^T A) A^T b$$

for appropriate function g (depending on λ).

Comparing ridge regression and PCR

- ▶ $\hat{w} = g(A^T A)A^T b$
- ▶ Ridge regression (with parameter λ): $g(z) = \frac{1}{z+\lambda}$
- ▶ PCR (with parameter λ): $g(z) = \mathbf{1}_{\{z \geq \lambda\}} \cdot \frac{1}{z}$
- ▶ Interpretation:

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

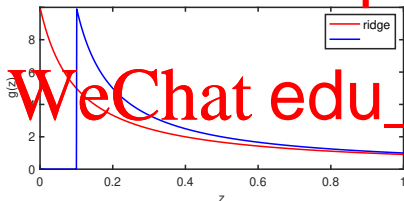


Figure 7: Comparison of ridge regression and PCR

- ▶ Let $A = \begin{bmatrix} \leftarrow & x_1^T & \rightarrow \\ & \vdots & \\ \leftarrow & x_n^T & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times d}$ (forget the $1/\sqrt{n}$ scaling)
- ▶ Try to approximate A with BC , where $B \in \mathbb{R}^{n \times k}$ and

<https://eduassistpro.github.io> high
utility

- ▶ Think of B as the encodings of the data
- ▶ “Dimension reduction” when $k < d$
- ▶ **Theorem** (Schmidt, 1907; Eckart-Young) Add WeChat edu_assist_pro
The best rank- k solution is given by truncating the singular value decomposition (SVD) of A

Singular value decomposition

- ▶ Every matrix $A \in \mathbb{R}^{n \times d}$ —say, with rank r —can be written as

$$A = \sum_{i=1}^r \sigma_i u_i u_i^\top$$

<https://eduassistpro.github.io>

- ▶ $v_1, \dots, v_r \in \mathbb{R}^d$ (orthonormal [ri](#))
- ▶ Can also write as

[Add WeChat edu_assist_pro](#)

where

- ▶ $U = [u_1 | \dots | u_r] \in \mathbb{R}^{n \times r}$, satisfies $U^\top U = I$
- ▶ $S = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$
- ▶ $V = [v_1 | \dots | v_r] \in \mathbb{R}^{d \times r}$, satisfies $V^\top V = I$

Truncated SVD

- ▶ Let A have SVD $A = \sum_{i=1}^r \sigma_i u_i v_i^\top$ (rank of A is r)
- ▶ Truncate at rank k (for any $k \leq r$): [rank- \$k\$ SVD](#)

Assignment Project Exam Help

$$A_k := \sum_{i=1}^k \sigma_i u_i v_i^\top$$

- ▶ <https://eduassistpro.github.io>

- ▶ $S_k = \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{k \times k}$

- ▶ $V_k = [v_1 | \dots | v_k] \in \mathbb{R}^{d \times k}$, satisfies

- ▶ **Theorem** (Schmidt/Eckart-Young)

$$\|A - A_k\|_F^2 = \min_{M: \text{rank}(M)=k} \|A - M\|_F^2 = \sum_{i=k+1}^r \sigma_i^2$$

Encoder/decoder interpretation (1)

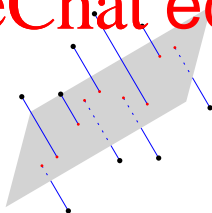
- ▶ Encoder: $x \mapsto \varphi(x) = V_k^T x \in \mathbb{R}^k$
 - ▶ Encoding rows of A : $AV_k = U_k S_k$
- ▶ Decoder: $z \mapsto V_k z \in \mathbb{R}^d$
 - ▶ Decoding rows of $U_k S_k$: $U_k S_k V_k^T = A_k$
- ▶ Same as k -dimensional PCA mapping!

Assignment Project Exam Help

<https://eduassistpro.github.io>

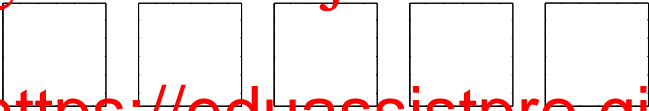
squared distances to data points.

Add WeChat edu_assist_pr



- ▶ Example: OCR data, compare original image to decoding of k -dimensional PCA encoding ($k \in \{1, 10, 50, 200\}$)

Assignment Project Exam Help


<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Application: Topic modeling (1)

- ▶ Start with n documents, represent using “bag-of-words” count vectors

- ▶ Arrange in matrix $A \in \mathbb{R}^{n \times d}$, where d is vocabulary size

| | aardvark | abacus | abalone |
|---|----------|--------|---------|
| ⋮ | | | |
| ⋮ | | | |
| ⋮ | | | |

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Application: Topic modeling (2)

- ▶ Rank k SVD provides an approximate factorization

Assignment Project Exam Help

where $B \in \mathbb{R}^{n \times k}$ and $C \in \mathbb{R}^{k \times d}$



- ▶ If rows of C were probability distrib

$C_{t,w}$ as probability that word w a

Add WeChat edu_assist_pr

Application: Matrix completion (1)

- ▶ Start with ratings of movies given by users
- ▶ Arrange in a matrix $A \in \mathbb{R}^{n \times d}$, where $A_{i,j}$ is rating given by user i for movie j .
- ▶ Netflix: $n = 480000$, $d = 18000$; on average, each user rates

- ▶ <https://eduassistpro.github.io>

$$B = \begin{bmatrix} \leftarrow b_1^\top \rightarrow \\ \vdots \\ \leftarrow b_n^\top \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C \in \mathbb{R}^{k \times d}$$

with goal of minimizing $\|A - BC\|_F^2$

- ▶ Note: If all entries of A were observed, we could do this with truncated SVD.

Application: Matrix completion (2)

- Need to find a low-rank approximation without all of A :

(low-rank) matrix completion

- Lots of ways to do this
- Popular way (used in Netflix competition): based on “stochastic gradient descent” (discussed later)

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Feature representations from matrix completion

- ▶ MovieLens data set ($n = 6040$ users, $d = 3952$ movies, $|\Omega| = 800000$ ratings)
- ▶ Fit B and C by using a standard matrix completion method (based on SGD, discussed later)
- ▶ k

<https://eduassistpro.github.io>

- ▶ Some nearest-neighbor pairs (c_j ,
 - ▶ Toy Story (1995), Toy Story 2 (1999)
 - ▶ Sense and Sensibility (1995), Emma (1996)
 - ▶ Heat (1995), Carlito's Way (1993)
 - ▶ The Crow (1994), Blade (1998)
 - ▶ Forrest Gump (1994), Dances with Wolves (1990)
 - ▶ Mrs. Doubtfire (1993), The Bodyguard (1992)