



eigenvalues. Eigenvectors for distinct eigenvalues are linearly independent.

- If there exists a non-trivial solution for  $q$ ,  $(A - \lambda I)$  is not invertible and characteristic polynomial  $\det(A - \lambda I) = 0$
- Eigendecomposition resp. diagonalization:**  $A = Q\Lambda Q^{-1}$  where  $Q$  is a matrix with the eigenvectors as columns and  $\Lambda$  is a diagonal matrix with the eigenvalues on the diagonal

- $\det(A) = \det(Q\Lambda Q^{-1}) = \prod_{i=1}^n \lambda_i$
- Symmetric eigendecomposition resp. unitary diagonalization:** For symmetric  $A$ :  $A = Q\Lambda Q^\top$  where  $Q$  is an orthogonal matrix with the eigenvectors as columns and  $\Lambda$  is a diagonal matrix with the eigenvalues on the diagonal
- Spectral theorem:** Square matrix  $A$  is symmetrically diagonalizable, iff  $AA^\top = A^\top A$
- Spectral theorem for symmetric matrices:** Every symmetric matrix  $A$  is symmetrically diagonalizable (due to Spectral theorem) and all its eigenvalues are real

**Positive definite (pd) and positive semi-definite matrices (psd)** —

- $A > 0$  iff  $x^\top A x > 0$
- $A \geq 0$  iff  $x^\top A x \geq 0$

Properties:

- If  $A$  is p(s)d,  $\alpha A$  is also p(s)d
- If  $A$  and  $B$  are p(s)d,  $A + B$  is also p(s)d
- If  $\det(A) = \prod_{i=1}^n \lambda_i > (\geq) 0$  resp.

$\{\lambda_i\}_{i=1}^n > (\geq) 0$  for pd (psd)

Pd properties:

- $I$  is pd
- If  $A$  is pd,  $A^{-1}$  is pd
- Cholesky decomposition:** If  $A$  is pd,  $A = BB^\top$
- If  $A$  and  $B$  are pd,  $(AB)^{-1} = B^{-1}A^{-1}$

Psd properties:

- If  $A$  is psd,  $BAB^\top$  is psd

## 2 Calculus

### Derivatives

**Rules** —

- Sum rule:  $\frac{\partial f+g}{\partial x} = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$
- Product rule:  $\frac{\partial f \times g}{\partial x} = f \times \frac{\partial g}{\partial x} + g \times \frac{\partial f}{\partial x}$
- Chain rule:  $\frac{\partial f(g)}{\partial x} = \frac{\partial f}{\partial g} \times \frac{g}{\partial x}$

**Common derivatives** —

- $\frac{\partial x^n}{\partial x} = nx^{n-1}$
- $\frac{\partial \log(x)}{\partial x} = \frac{1}{x}$
- $\frac{\partial e^{kx}}{\partial x} = k \times e^{kx}$
- $\frac{\partial \sqrt{x}}{\partial x} = \frac{1}{2\sqrt{x}}$

**Partial and directional derivative** —

- For a function that depends on  $n$  variables  $\{x_i\}_{i=2}^n$ , partial derivative is slope of tangent line along direction of one specific variable  $x_i$
- Directional derivative is slope of tangent line along direction of selected unit vector  $u$

**Gradient** —

- Given scalar-valued function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , returns vector containing first-order partial derivatives:

$$\nabla_x f : [\frac{\partial f}{\partial x_1} \dots \frac{\partial f}{\partial x_n}]^\top$$

- Gradient points in direction of greatest upward slope of  $f$
- Magnitude of gradient equals rate of change when moving into direction of greatest upward slope

**Hessian** —

- Given scalar-valued function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , returns matrix containing second-order partial derivatives:

$$\mathcal{H} = \nabla_x^2 f : \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

- $\mathcal{H}$  is symmetric
- Jacobian** —
- Given vector-valued function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  with  $f = [f_1(x), \dots, f_m(x)]^\top$ , returns matrix containing first-order partial derivatives:

$$\nabla_x f : \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

**Matrix calculus rules** —

- $\frac{\partial a^\top x}{\partial x} = a$
- For square  $A$ :  $\frac{\partial x^\top A x}{\partial x} = (A + A^\top)x$
- $\frac{\partial a^\top A b}{\partial A} = ab^\top$
- For symmetric  $A$ :  $\frac{\partial x^\top A x}{\partial x} = 2Ax$
- $\frac{\partial a^\top A^{-1} b}{\partial A} = -(A^\top)^{-1} a c^\top (A^\top)^{-1}$
- $\frac{\partial \log(|A|)}{\partial A} = (A^\top)^{-1}$

### Extrema

**Conditions for local minima and maxima** —

- Point is a stationary point, i.e. first-order derivative = 0
- If Hessian is pd, it's a local minimum, if Hessian is nd, it's a local maximum, if Hessian is indefinite, it's a saddle point
- Local minima and maxima are the unique global minima and maxima in strictly convex functions resp. one of possibly infinitely many global minima and maxima in convex functions

**Convexity** —

- For a convex function:
  - $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$
  - Hessian of stationary point(s) is psd
  - Global minimum exists, but may not be unique
- For a strictly convex function:
  - $f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$
  - Hessian of stationary point is pd
  - Unique global minimum exists
- Sum of convex functions is also convex
- Sum of convex and strictly convex function is strictly convex

**Nature of optimum** — What does Hessian and function look like?

- If Hessian is pd and loss function is strictly convex, stationary point is a global minimum, and there is a unique solution
- If Hessian is psd and loss function is convex, stationary point is a global minimum, and there may be a geometrically unique or infinitely many solutions
- If Hessian is p(s)d but loss function is not convex, stationary point may be a local

minimum and there may be a geometrically unique or infinitely many solutions

**Optimization approach** — Is function differentiable, continuous, and are relevant terms invertible?

- If yes, analytically solvable
- If no, numerically solvable (e.g. via gradient descent)

**Constrained optimization** —

- Lagrangian function:  $\mathcal{L}(x, \lambda) = f(x) + \lambda g(x)$ , where  $g(x)$  is an  $(m - 1)$  dimensional constraint surface and  $\lambda$  is the Lagrange multiplier
- $\nabla_x \mathcal{L} = \nabla_x f(x) + \lambda \nabla_x g(x)$
- $\nabla_\lambda \mathcal{L} = g(x)$

For equality constraints: Minimize  $f(x)$  subject to  $g(x) = 0$

- Gradient of  $f(x)$  must be orthogonal to constraint surface, otherwise (if it points into any direction along the constraint surface)  $f(x)$  could still decrease for movements along the constraint surface
- On the constraint surface,  $g(x)$  is a constant, so moving along any direction on the constraint surface has a directional derivative of 0. Since the gradient of  $g(x)$  points into the direction of steepest ascent, it must be orthogonal to the constraint surface, otherwise (if it points into any direction along the constraint surface)  $g(x)$  would not be constant on the constraint surface
- Then, gradients are parallel at optimum:  $\nabla_x f(x^*) = \lambda \times \nabla_x g(x^*)$
- To find  $x^*$  and  $\lambda^*$ :
  - $\nabla_x L = 0$ , expresses parallelity condition at minimum  $x^*$
  - $\nabla_\lambda L = 0$ , expresses constraint
  - This is an unconstrained optimization problem
- Optimum  $x^*$  and  $\lambda^*$  represents a saddle point of  $\mathcal{L}$

For inequality constraints: Minimize  $f(x)$  subject to  $g(x) \leq 0$

- If  $x^*$  lies in  $g(x) < 0$ , constraint is inactive
- Otherwise, if  $x^*$  lies in  $g(x) = 0$ , constraint is active:
  - Gradient of  $f(x)$  must point towards  $g(x) < 0$  region, otherwise (if it would point away from  $g(x) < 0$  region) the optimum would lie in this region
  - Then, gradients are anti-parallel at optimum:  $\nabla_x f(x^*) = -\lambda \times \nabla_x g(x^*)$
- To find  $x^*$  and  $\lambda^*$ :
  - $\nabla_x L = 0$  subject to **Karush Kuhn Tucker conditions**:
    - $g(x) \leq 0$
    - $\lambda \geq 0$
    - Complementary slackness condition**:  $\lambda g(x) = 0$ , with  $\lambda = 0, g(x) < 0$  for inactive constraints and  $\lambda > 0, g(x) = 0$  for active constraints
  - $\nabla_\lambda \mathcal{L} = 0$  given complementary slackness condition
  - This is not an unconstrained optimization problem, but can be solved via duality

- Optimum  $x^*$  and  $\lambda^*$  represents a saddle point of  $\mathcal{L}$

For multiple constraints: Minimize  $f(x)$  subject to  $m$  inequality constraints

$$\{g^{(i)}(x) \leq 0\}_{i=1}^m \text{ and } p \text{ equality constraints } \{h^{(j)}(x) = 0\}_{j=1}^p$$

- Then, Lagrangian is given by:  $\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_{i=1}^m \mu^{(i)} g^{(i)}(x) + \sum_{j=1}^p \lambda^{(j)} h^{(j)}(x)$
- Then, general solution  $x^*, \lambda^*, \mu^*$  is given by:  $\nabla_x \mathcal{L} = 0$  subject to:
  - $\{g^{(i)}(x) \leq 0\}_{i=1}^m$  and  $\{h^{(j)}(x) = 0\}_{j=1}^p$
  - $\{\mu^{(i)} \geq 0\}_{i=1}^m$
  - $\{\mu^{(i)} g^{(i)}(x) = 0\}_{i=1}^m$

Solving inequality constraints via duality – **primal problem**:

- $\min_x [\max_{\lambda, \mu} \mathcal{L}]$

- $\max_{\lambda, \mu} \mathcal{L} = f(x) + \max_{\lambda, \mu} [\sum_{i=1}^m \mu^{(i)} g^{(i)}(x) + \sum_{j=1}^p \lambda^{(j)} h^{(j)}(x)]$

- Second term gives rise to barrier function:
  - $= 0$  subject to constraints being met, given complementary slackness condition for inequality constraints and  $h^{(j)}(x) = 0$  for equality constraints, which implies that primal problem becomes  $\min_x (f(x))$
  - $= \infty$  otherwise, which implies that primal problem cannot be solved

Solving inequality constraints via duality – **weak duality**:

- Given minimax theorem,  $\min_x [\max_{\lambda, \mu} \mathcal{L}] \geq \max_{\lambda, \mu} [\min_x \mathcal{L}]$ , which gives a lower bound of minimum of primal problem
- $\min_x \mathcal{L}$  is an unconstrained optimization problem
- $\max_{\lambda, \mu} [\min_x \mathcal{L}]$  is a concave maximization problem

Solving inequality constraints via duality – **strong duality**:

- If constraint qualifications are fulfilled,  $\min_x [\max_{\lambda, \mu} \mathcal{L}] = \max_{\lambda, \mu} [\min_x \mathcal{L}]$
- $\min_x \mathcal{L}$  can be solved for general solution  $x^*$  in terms of  $\lambda, \mu$
- Plug  $x^*$  back into  $\mathcal{L}$  and maximize to find solutions  $\lambda^*, \mu^*$
- Specify  $x^*$  based on  $\lambda^*, \mu^*$

## 3 Probability and Statistics

### Terminology

**Kolmogorov axioms** — Probability space defined by:

- Sample space: All possible outcomes  $\Omega = \{\omega_1, \dots, \omega_n\}$
  - Event space: All possible results, where an event is a subset of the sample space
  - Probability measure: Function that assigns a probability to an event
- Axioms:
- Event space must be a **sigma algebra**:
    - If  $A$  is in sample space, its complement is also in sample space
    - If  $A_1, \dots A_n$  are in sample space, their union is also in sample space

- Probability measure must satisfy:
  - $0 \leq \mathbb{P}(A) \leq 1$
  - $\mathbb{P}(\Omega) = 1$
  - If  $A_1, A_2, \dots$  are in sample space and do not intersect, then  $\mathbb{P}(A_1 \cup A_2 \cup \dots) = \int_{n=1}^\infty \mathbb{P}(A_n)$

Further properties:

- All sets than can be formed from left and right inclusive interval  $[0, a]$  are events.

On that basis:  $(b, 1] = [0, b]^c \in$  event space.

**Variables** —

- Random variable:
  - Discrete random variable: Characterized by pmf
  - Continuous random variable: Characterized by pdf
- Independent random variables:
  - $\mathbb{P}(A|B) = \mathbb{P}(A)$  and  $\mathbb{P}(B|A) = \mathbb{P}(B)$
  - $\mathbb{E}(AB) = \mathbb{E}(A)\mathbb{E}(B)$
  - Correlation is 0
  - $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$
  - Functions of independent random variables are also independent
- Conditionally independent random variables: Two random variables  $\mathcal{X}$  and  $\mathcal{Y}$  are conditionally independent, if there is a confounder  $\mathcal{L}$  that causally affects both variables, but if we control for this confounder, the variables are not causally connected
- I.I.D. random variables: Independent and from identical distribution

**Events** —

- Complement:  $\mathbb{P}(A^C) = 1 - \mathbb{P}(A)$  and  $\mathbb{P}(A \cup A^C) = \mathbb{P}(A)\mathbb{P}(A^C)$
- Disjoint / mutually exclusive vs. joint / mutually inclusive
- Subset  $A \subset B$  with  $\mathbb{P}(A) < \mathbb{P}(B)$

**Probabilities** —

- Marginal probability  $\mathbb{P}(A)$ : Probability for single variable:  $p(\mathcal{X}) = \sum_{\mathcal{Y}} p(x, y)$  resp.

$$f(\mathcal{X}) = \int_{\mathcal{Y}} f(x, y) dy$$

- Joint probability  $\mathbb{P}(A \cap B)$ : Probability for combination of variables, given by all possible combinations resp. convolution of their pdfs

- Conditional probability  $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ : Probability for variable, given other variable:  $p(\mathcal{X}|\mathcal{Y}) = \frac{p(x, y)}{\sum_{\mathcal{X}} p(x, y)}$  resp.

$$f(\mathcal{X}|\mathcal{Y}) = \frac{f(x, y)}{\int_{\mathcal{X}} f(x, y) dy}$$

- $\mathbb{P}(A|B) = 1 - \mathbb{P}(A^C|B)$
- $\mathbb{P}(A_1|B) + \mathbb{P}(A_2|B) + \dots = 1$

- Bayesian terminology**:
  - Prior  $\mathbb{P}(\text{parameter})$
  - Posterior  $\mathbb{P}(\text{parameter}|\text{data})$
  - Likelihood  $\mathbb{P}(\text{data}|\text{parameter})$
  - Evidence  $\mathbb{P}(\text{data})$
- Bayes theorem**: Posterior  $\mathbb{P}(A|B) = \frac{\text{Likelihood } \mathbb{P}(B|A) \times \text{Prior } \mathbb{P}(A)}{\text{Evidence } \mathbb{P}(B)}$

### Measures

**Expected value** —  $\mathbb{E}(\mathcal{X}) = \sum_{\mathcal{X}} x \times p(x)$  resp.

$\mathbb{E}(\mathcal{X}) = \int_{-\infty}^\infty x \times f(x) dx$  with pmf resp. pdf — Properties:



- $\mathbb{E}(\alpha) = \alpha$
- $\mathbb{E}(\alpha\mathcal{X} + \beta) = \alpha\mathbb{E}(\mathcal{X}) + \beta$
- $\mathbb{E}(\alpha\mathcal{X} + \beta\mathcal{Y}) = \alpha\mathbb{E}(\mathcal{X}) + \beta\mathbb{E}(\mathcal{Y})$
- For orthogonal variables:  
 $\mathbb{E}((\mathcal{X} + \mathcal{Y})^2) =$

*Cauchy Schwarz inequality:*

$$\mathbb{E}(\mathcal{X}, \mathcal{Y})^2 \leq \mathbb{E}(\mathcal{X}^2)\mathbb{E}(\mathcal{Y}^2)$$

*Standard deviation* —  $\sqrt{\mathbb{V}(\mathcal{X})}$

*Covariance* —

- Univariate variance of a random variable:  
 $\mathbb{V}(\mathcal{X}) = \mathbb{E}((\mathcal{X} - \mathbb{E}(\mathcal{X}))^2) = \mathbb{E}(\mathcal{X}^2) - \mathbb{E}(\mathcal{X})^2$   
 where  $\mathbb{E}(\mathcal{X}^2)$  is the unnormalized correlation resp. inner product
- Univariate covariance of two random variables:  $\text{Cov}(\mathcal{X}, \mathcal{Y}) = \mathbb{E}((\mathcal{X} - \mathbb{E}(\mathcal{X}))(\mathcal{Y} - \mathbb{E}(\mathcal{Y}))) = \mathbb{E}(\mathcal{X}\mathcal{Y}) - \mu_{\mathcal{X}}\mu_{\mathcal{Y}}$   
 where  $\mathbb{E}(\mathcal{X}\mathcal{Y})$  is the unnormalized correlation resp. inner product
- Multivariate covariance matrix of a vector:

$$\begin{aligned} - \Sigma &= \text{Cov}(\mathcal{X}) = \mathbb{E}((\mathcal{X} - \mathbb{E}(\mathcal{X}))(\mathcal{X} - \mathbb{E}(\mathcal{X}))^\top) = \mathbb{E}(\mathcal{X}\mathcal{X}^\top) - \mathbb{E}(\mathcal{X})\mathbb{E}(\mathcal{X})^\top = \\ &\begin{bmatrix} \text{Var}(\mathcal{X}_1) & \dots & \text{Cov}(\mathcal{X}_1, \mathcal{X}_m) \\ \vdots & \ddots & \vdots \\ \text{Cov}(\mathcal{X}_m, \mathcal{X}_1) & \dots & \text{Var}(\mathcal{X}_m) \end{bmatrix} \\ &\text{where } R = \mathbb{E}(\mathcal{X}\mathcal{X}^\top) \text{ is the unnormalized correlation matrix} \\ - \Sigma &\text{ and } R \text{ are symmetric and psd} \\ - \Sigma &= R - \mu_{\mathcal{X}}\mu_{\mathcal{X}}^\top \end{aligned}$$

Properties - variance:

- $\mathbb{V}(\alpha) = 0$
- $\mathbb{V}(\alpha\mathcal{X} + \beta) = \alpha^2\mathbb{V}(\mathcal{X})$
- $\mathbb{V}(\mathcal{X} + \mathcal{Y}) = \mathbb{V}(\mathcal{X}) + 2\text{Cov}(\mathcal{X}, \mathcal{Y}) + \mathbb{V}(\mathcal{Y})$
- For uncorrelated (and independent) variables:
- For vector  $y = Ax$ :  
 $\mathbb{V}_y = A\mathbb{V}_X A^\top$

Properties - covariance:

- $\text{Cov}(\mathcal{X}, \mathcal{X}) = \mathbb{V}(\mathcal{X})$
- $\text{Cov}((\alpha\mathcal{X} + \beta\mathcal{Y}), \mathcal{Z}) = \alpha\text{Cov}(\mathcal{X}, \mathcal{Z}) + \beta\text{Cov}(\mathcal{Y}, \mathcal{Z})$
- If covariance of two random variables is 0, they are uncorrelated, but not necessarily independent. Then,  $\mathbb{E}(\mathcal{X}\mathcal{Y}) = \mathbb{E}(\mathcal{X})\mathbb{E}(\mathcal{Y})$
- If covariance and unnormalized correlation of two random variables is 0, they are orthogonal, but not necessarily independent. Then,  $\mathbb{E}(\mathcal{X}\mathcal{Y}) = 0$
- For vector  $y = Ax$ :  
 $-\Sigma_y = A\Sigma_X A^\top$   
 $-R_y = AR_X A^\top$

*Cauchy Schwarz inequality:*

- $\text{Cov}(\mathcal{X}, \mathcal{Y})^2 \leq \mathbb{V}(\mathcal{X})\mathbb{V}(\mathcal{Y})$
- $\mathbb{E}(\mathcal{X}\mathcal{Y})^2 \leq \mathbb{E}(\mathcal{X}^2)\mathbb{E}(\mathcal{Y}^2)$

*Correlation* — Normalized covariance

- Univariate correlation of a random variable:  $\text{Cor}(\mathcal{X}, \mathcal{Y}) = \frac{\text{Cov}(\mathcal{X}, \mathcal{Y})}{\sqrt{\mathbb{V}(\mathcal{X})}\sqrt{\mathbb{V}(\mathcal{Y})}}$

- Multivariate correlation matrix of a vector:  
 $-P = \text{Cor}(\mathcal{X}) =$

$$\begin{bmatrix} 1 & \dots & \text{Cor}(\mathcal{X}_1, \mathcal{X}_m) \\ \vdots & \ddots & \vdots \\ \text{Cor}(\mathcal{X}_m, \mathcal{X}_1) & \dots & 1 \end{bmatrix}$$

- $-P$  is symmetric and psd
- Correlation is bounded between 0 and 1, given Cauchy Schwarz Inequality
- If correlation of two random variables is 0, they are not necessarily independent

**Probability Distributions**

PMF, CDF, PDF —

- Cumulative density function  $F(r)$  (CDF):  
 $F(r) = p(x \leq r)$
- Probability mass function  $p(x)$  (PMF) for discrete random variables:  $p(x)$
- Probability density function (PDF)  $f(x)$  for continuous random variables:  
 $\int_{-\infty}^r f(x)dx = p(x \leq r) = F(r)$
- Properties of CDF and PDF:
  - Derivative of CDF returns PDF, integral of PDF returns CDF
  - Monotonically non-decreasing: If  $s < r$ ,  $F(s) < F(r)$
  - $\lim_{r \rightarrow -\infty} F(r) = 0$
  - $\lim_{r \rightarrow \infty} F(r) = 1$
  - Right-continuous:  $\lim_{s \rightarrow r^+} F(s) = F(r)$
  - $\lim_{s \rightarrow r^-} F(s) = F(x < r) = F(s) - F(x = r)$
  - $\int_a^b f(x)dx = F(b) - F(a) = p(a < x \leq b)$
  - $\int_{-\infty}^{\infty} f(x)dx = 1$

*Normal distribution* —  $\mathcal{X} \sim \mathcal{N}(\mu, \sigma^2)$

For univariate, PDF:  $\frac{1}{\sqrt{2\pi\sigma}} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$

For multivariate, PDF:  
 $\frac{1}{2\pi\sigma^{n/2}} \frac{1}{|\Sigma|^{1/2}} \exp(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu))$  where the term in the exponent is a quadratic form  
*Bernoulli distribution* — trial with success (probability  $p$ ) or failure (probability  $1-p$ )

- PDF:  $p(x)p^x(1-p)^{1-x}$
- Mean:  $\mathbb{E}(x) = p$
- Variance:  $\mathbb{V}(x) = p(1-p)$

*Binomial distribution* —  $n$  independent Bernoulli trials with  $k$  successes

- PDF:  $\binom{n}{k} p^k (1-p)^{n-k}$
- Mean:  $\mathbb{E}(x) = np$
- Variance:  $\mathbb{V}(x) = np(1-p)$

*Poisson distribution* —

- PDF:  $e^{-\lambda} \frac{\lambda^x}{x!}$
- Mean:  $\mathbb{E}(x) = \lambda$
- Variance:  $\mathbb{V}(x) = \lambda$

*Beta distribution* —

- PDF:  $\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$
- Mean:  $\mathbb{E}(x) = \frac{\alpha}{\alpha+\beta}$
- Variance:  $\mathbb{V}(x) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

**Laws of Large Numbers and Inequalities**

- Laws of large numbers* — Sample mean of iid variables converges to population mean as  $n \rightarrow \infty$
- Jensen's inequality* — Relates expected value of a convex function of a random variable to the convex function of the expected value of that random variable  
 $\mathbb{E}(f(\mathcal{X})) \geq f(\mathbb{E}(\mathcal{X}))$
- Markov's inequality* —  $p(x \geq t) \leq \frac{\mathbb{E}(x)}{t}$   
 Interesting only for  $t \geq \mathbb{E}(x)$  because  $p(x \geq t)$  must then be less than or equal to 1

Generalizations:

- $p(|x| \geq t) \leq \frac{\mathbb{E}(|x|)}{t}$
- $p(|x| \geq t) \leq \frac{\mathbb{E}(|x|^n)}{t^n}$

*Chebyshev's inequality* —

$$p(|x - \mu_x| \geq \alpha|\sigma_x|) \leq \frac{1}{\alpha^2}$$

Interesting only for  $\alpha > 1$   
 Implications:

- For  $n$  variables:  $p(|S_n - \mu_x| \geq \epsilon) \leq \frac{\sigma_x^2}{n\epsilon^2}$  where  $S_n$  is the sample mean

**4 Information Theory**

**Description**

*Entropy* —

- $H(x) = -\sum_x p(x) \log(p(x)) = -\sum_{x,y} p(x,y) \log(p(x))$  resp.  
 $H(x) = -\int p(x) \log(p(x)) dx$
- Measure of randomness in a variable resp. quantifies uncertainty of a distribution
- Properties:
  - $H(x) \geq 0$
  - $H(x)$  is maximized, when  $x$  is a uniform random variable
  - For independent variables:  
 $H(x,y) = H(x) + H(y)$

*Conditional entropy* —

- $H(x|y) = -\sum_{x,y} p(y) p(x|y) \log(p(x|y)) = -\sum_{x,y} p(x,y) \log(\frac{p(x,y)}{p(y)})$
- Measure of how much information of  $x$  is revealed by  $y$

Properties:

- $0 \leq H(x|y) \leq H(x)$  with equality if when  $x$  is independent with  $y$  resp. if  $y$  completely determines  $x$

*Mutual information* —

- $I(x;y) = H(x) - H(x|y) = -\sum_{x,y} p(x,y) \log(\frac{p(x)p(y)}{p(x,y)})$
- Measure of how much information of  $x$  is left after  $y$  is revealed

Properties:

- $0 \leq I(x;y) \leq H(x)$  with equality if  $y$  completely determines  $x$  resp. if  $x$  is independent with  $y$

*KL divergence* —

- $KL(p;q) = \sum_x p(x) \log(\frac{p(x)}{q(x)})$
- Measures the extra information or inefficiency when approximating a true distribution over  $x$ ,  $p$ , with a predicted one,  $q$

Properties:

- $KL(p;q) \geq 0$
- Cross entropy* —  
 $CE(p|q) = KL(p;q) + H(p) = -\sum_x p(x) \log(q(x))$

- Measures the total uncertainty when using the predicted distribution  $q$  to represent the true distribution  $p$ , combining both the model's error and the intrinsic uncertainty of the true distribution

Properties:

- $KL(p;q) \geq 0$

**5 ML Paradigms**

**Frequentism**

*Description* —

- Parametric approach

- $\theta$  as fixed, unknown quantity,  $X$  as random, and known quantity
- Makes point estimate
- Focuses on maximizing likelihood  $p(X|\theta)$  to infer posterior  $p(\theta|X)$
- Only requires differentiation methods
- High variance, but low bias

*MLE estimator*

- Maximizes log-likelihood:  $\hat{\theta} = \text{argmax}_{\theta} (L) = \text{argmax}_{\theta} (\prod_{i=1}^n p(x_i|\theta)) = \text{argmax}_{\theta} (\sum_{i=1}^n \log(p(x_i|\theta)))$

- Advantages
  - Consistent:  $\hat{\theta} \rightarrow \theta$  as  $n \rightarrow \infty$
  - Asymptotically normal:  $\frac{1}{\sqrt{n}}(\hat{\theta} - \theta)$  covers to  $\mathcal{N}(0, J^{-1}(\theta)I(\theta)J^{-1}(\theta))$

where  $J = -\mathbb{E}[\frac{\partial^2 \log(p(x|\theta))}{\partial \theta \partial \theta^\top}]$  and where  $I$  is the Fisher information

- Asymptotically efficient:  $\hat{\theta}$  minimizes  $\mathbb{E}[(\hat{\theta} - \theta)^2]$  as  $n \rightarrow \infty$ 
  - Not necessarily the best estimator, especially for small samples in a multivariate context
  - (cf. Rao-Cramer bound)
- Equivariant: If  $\hat{\theta}$  is MLE of  $\theta$ , then  $g(\hat{\theta})$  is MLE of  $g(\theta)$

- Proofs of advantages
  - Asymptotically normal:
    - We start with the score and set it to 0 for optimization with regard to  $\theta$ :  
 $\Lambda = \frac{\partial}{\partial \theta} \log(p(x|\theta)) = 0$
    - With a Taylor expansion, we can show that  $(\hat{\theta} - \theta)\sqrt{n} = \frac{1}{\sqrt{n}} \Lambda [-\frac{1}{n} \frac{\partial^2}{\partial \theta \partial \theta^\top} \sum_{i=1}^n \log(p(x_i|\theta))]^{-1}$  where  $\Lambda$  is the score
    - We set  
 $J = [-\frac{1}{n} \frac{\partial^2}{\partial \theta \partial \theta^\top} \sum_{i=1}^n \log(p(x_i|\theta))]$
    - $\frac{1}{\sqrt{n}} \Lambda$  is a random vector with covariance matrix  $I$  and converges to the normal distribution  $\sim \mathcal{N}(0, I)$
    - Then,  
 $(\hat{\theta} - \theta)\sqrt{n} = J^{-1} \frac{1}{\sqrt{n}} \Lambda \sim J^{-1} \mathcal{N}(0, I)$
    - $\mathbb{V}(J^{-1} \frac{1}{\sqrt{n}} \Lambda) = \mathbb{E}[J^{-1} I J^{-1}]$
    - This equality is given because  $\mathbb{V}(x) = \mathbb{E}[x - \mathbb{E}(x)] = \mathbb{E}[x]$  if  $\mathbb{E}(x) = 0$ , which is the case here, given that the expected score is 0
    - So we have shown that  
 $\hat{\theta} - \theta \sqrt{n} = J^{-1} \frac{1}{\sqrt{n}} \Lambda \sim \mathcal{N}(0, J^{-1} I J^{-1})$

- Equivariant:
  - Let  $t = g(\theta)$  and  $h = g^{-1}$
  - Then,  $\theta = h(t) = h(g(\theta))$
  - For all  $t$  we have:  $L(t) = \prod_i p(x^{(i)}|h(t)) = p(x^{(i)}|\theta) = L(\theta)$
  - Hence, for all  $t$  we can say:  
 $L(t) = L(\theta)$  and  $L(\hat{t}) = L(\hat{\theta})$

*PAC estimator*

- Generates probabilistic bounds for parameter  $\theta$  that is approximately known with a high probability:
  - Probability of being correct:  $1 - \delta$

- Degree of approximation:  $\epsilon$
- Given Hoeffding's inequality, the probability that the error is greater than  $\epsilon$  is bounded

**Bayesianism**

*Description* —

- Parametric approach
- $\theta$  as random, unknown quantity,  $X$  as random, and known quantity
- Makes estimate in form of distribution
- Focuses on leveraging prior and likelihood to infer posterior:

$$p(\theta|X, y) = \frac{p(\theta)p(y|X, \theta)}{p(y|X)} = \frac{p(\theta)p(y|X, \theta)}{\int p(\theta)p(y|X, \theta)d\theta} \propto p(\theta)p(y|X, \theta)$$

- Requires integration methods for normalizing constant in denominator, which can be intractable, in which case MAP estimator can provide an alternative
- Low variance, but high bias

*Mean estimator*

- Takes expected value and variance posterior
- Returns estimate that reflects central tendency and overall uncertainty

*MAP estimator*

- Maximizes posterior (i.e. takes point where posterior density is highest):  
 $\hat{\theta} = \text{argmax}_{\theta} (p(\theta|X))$
- Returns single point estimate

**Statistical Learning**

*Description* —

- We want to minimize expected risk  $\mathcal{R}(f) = \mathbb{E}_{X,Y}[1f(X) \neq Y]$ , but this is difficult because
  - We don't have access to the joint distribution of  $X, Y$
  - We cannot find  $f$ , without any assumptions on its structure
  - It's unclear how to minimize the expected value
- Therefore, we make following choices:
  - We collect sample  $Z$
  - We restrict space of possible choices of  $f$  to a set  $\mathcal{H}$
  - We use a loss function to approximate the expected value
- With these choices, we approximate the expected risk via the empirical risk  
 $\hat{\mathcal{R}}(f) = \hat{L}(Z, f) = \frac{1}{n} \sum_i L(y_i, f(x_i))$

**6 Model Taxonomy**

**Supervised vs. Unsupervised Learning**

**TBA**

**Active Learning**

- Assume:
  - Domain space  $\mathcal{X}$
  - Sample space  $S \subseteq \mathcal{X}$
  - Labeled data  $D_{n-1}(x_i, y_i)_{i < n}$
  - Target space  $\mathcal{A} \subseteq \mathcal{X}$
  - We estimate  $y_x = f_x + \epsilon_x$
- We aim to find the next  $x_n$  that gives us the most information about  $f$  in  $\mathcal{A}$
- Information gain can be quantified as maximizing the conditional mutual information between  $y_x$  and  $f$ :

$$IG[f_x, y_x|D_{n-1}] = H(D_{n-1}) - H(D_{n-1}|x_n)$$

where  $H(D_{n-1})$  is the uncertainty about  $D$

before labeling  $x_n$  and  $H(D_{n-1}|x_n)$  is the uncertainty about  $D$  after labeling  $x_n$ . We want to minimize the latter, i.e. we want to maximize the delta between the former and the latter

- We pick  $x_n = \operatorname{argmax}_{x \in S} IG[f_x, y_x|D_{n-1}]$
- To find a closed-form solution, we assume that  $f$  is a Gaussian process with a known mean and kernel function:
  - $f \sim \mathcal{GP}(\mu, k)$
  - $f = (f_{x_1}, f_{x_2}, \dots) \sim \mathcal{N}(\mu, \Sigma)$  where elements in mean vector are  $\mu_i = \mu(x_i)$  and elements in covariance matrix are  $\Sigma_{ij} = k(x_i, x_j)$

- Under this assumption, we can show that 
$$IG[f_x, y_x|D_{n-1}] = \frac{1}{2} \log \left( \frac{\mathbb{V}(y_x|D_{n-1})}{\mathbb{V}(y_x|f_x, D_{n-1})} \right)$$

*Safe Bayesian learning* —

- Bayesian approach to active learning
- Assume:
  - We have stochastic process  $f^*$
  - We can iteratively choose points  $x_1, \dots, x_{n-1} \in \mathcal{X}$  and observe  $y_i = f^*(x_1), \dots, y_{n-1} = f^*(x_{n-1})$
  - Points should lie in safe area  $S^*$  which is the set of  $x \in \mathcal{X}$  such that another stochastic process  $g^*(x) \geq 0$
  - For chosen points, we can also observe  $z_i = g^*(x_1), \dots, z_{n-1} = g^*(x_{n-1})$  which are measurements of confidence, indicating high confidence when above 0
- We aim to find estimates of sample space  $S$  and target space  $\mathcal{A}$
- To do so, we fit a Gaussian process on observed  $\{(x_i, y_i)\}_{i < n}$  and  $\{(x_j, z_j)\}_{j < n}$ . Gaussian process over  $f$  and  $g$  induces two bounds respectively, which provide the 95% confidence interval of  $\mathbb{E}[f(x)]$  resp.  $\mathbb{E}[g(x)]$ :

- Upper bound function  $u_n^f(x)$  resp.  $u_n^g(x)$
- Lower bound function  $l_n^f(x)$  resp.  $l_n^g(x)$
- Gaussian process over  $g$  allows to derive pessimistic and optimistic estimate of safe area:
  - Pessimistic:  $S_n = \{x : l_n^g(x) \geq 0\}$
  - Optimistic:  $\hat{S}_n = \{x : u_n^g(x) \geq 0\}$
- We then gather estimates, where upper bound of  $f$  lies above baseline set by maximum value of lower bound of  $f$ :

$$\mathcal{A}_n = \{x \in \hat{S}_n : u_n^f(x) \geq \max_{x' \in S_n} l_n^f(x')\}$$

- We can then perform active learning with sample space  $S = S_n$  and target space  $\mathcal{A} = \mathcal{A}_n$

*Batch active learning* —

- Variant of active learning
- Assume:
  - Domain space  $\mathcal{X}$  and distribution  $P$  over  $\mathcal{X}$
  - Oracle to unknown function  $f : \mathcal{X} \rightarrow \mathcal{Y}$
  - Population set  $\mathcal{X} = \{x_1, \dots, x_m\} \subseteq \mathcal{X}$
  - Budget  $b \leq m$
- We aim to find next batch of data points  $L \subseteq \text{mathcal{X}}$  subject to  $|L| = b$  that gives us the most information
- Suppose we know  $Z = \{(x, f(x)) : x \in L\}$

- 1-nearest-neighbor classifier  $\hat{f}$  is fitted to  $Z$
- Let  $B_\delta(x) = \{x' \in \mathcal{X} : \|x - x'\| \leq \delta\}$  be the set of sufficiently close points to  $x$ 
  - We consider  $B_\delta(x)$  pure if  $f$  yields same results for all of  $B_\delta(x)$
- Impurity of  $\delta$  is given by  $\hat{\pi}(\delta) = P(\{x \in \mathcal{X} : B_\delta(x) \text{ is not pure}\})$
- Let  $C(L, S) = \bigcup_{x \in L} B_\delta(x)$  be the union of all sets  $B$ 
  - $C = C_r \cup C_w = \{x \in C : \hat{f}(x) = f(x)\} \cup \{x \in C : \hat{f}(x) \neq f(x)\}$
- We have  $C_w \subseteq \{x \in \mathcal{X} : B_\delta(x) \text{ is not pure}\}$ . Then,  $P(C_w) \leq \hat{\pi}(\delta)$
- $\{x : \hat{f}(x) \neq f(x)\} \subseteq C_w \cup C_r^C \subseteq \{x \in \mathcal{X} : B_\delta(x) \text{ is not pure}\} \cup C^C$
- Then, we have  $\mathcal{R}(f) = P(\hat{f}(x) \neq f(x) \leq \hat{\pi}(\delta) + 1 - P(C)$
- We need to choose  $L$  and  $\delta$  such that  $\mathcal{R}(\hat{f})$  is minimized
- We approach this by minimizing the upper bound, by picking  $\delta$  and choosing  $C$  that maximizes  $P(C)$ :  $\operatorname{argmax}_{L \subseteq \mathcal{X}, |L|=b} P(\bigcup_{x \in L} B_\delta(x))$
- Two challenges:
  - We don't know the distribution
  - Problem is NP-hard
- We address 1) by using the empirical distribution induced by  $X$ . Then, we have:  $\operatorname{argmax}_{L \subseteq \mathcal{X}, |L|=b} \frac{1}{|X|} \|\{x' : \|x' - x\| \leq \delta, \text{ for some } x \in L\}\|$
- We address 2) with greedy algorithm:
  - Input:  $x \subseteq \mathcal{X}, b \in \mathbb{N}$
  - Output:  $L \subseteq X$  of size  $b$ 
    - $G = (x, E)$  where  $E = \{(x, x') : \|x - x'\| \leq \delta\}$
    - $L = \emptyset$
    - For  $i = 1, \dots, b$ :
      - $\hat{x} \leftarrow \operatorname{argmax}_{x \in \mathcal{X}} |\{x' : (x, x') \in E, x \in \mathcal{X}\}|$
      - $L \leftarrow L \cup \hat{x}$
      - $E \leftarrow E - (\{\hat{x}\} \times (B_\delta(\hat{x}) \cap x))$
- Return  $L$

### 7 Model Optimization

#### Gradient Descent

Numeric optimization procedure

*Gradient descent* —

- Uses entire training set to evaluate whether new parameter is more optimal than previous one
- Slow and less likely to escape local minima due to randomness, but accurate
- Algorithm:
  - Set  $\eta > 0$
  - Randomly initialize  $\beta_{(t=0)}$
  - $\beta_{(t+1)} \leftarrow \beta_{(t)} - \eta \nabla_\beta LO|_{\beta=\beta_{(t)}}$
  - $t \leftarrow t + 1$
  - Repeat 3 and 4 until  $\nabla_\beta LO = 0$
- Stochastic gradient descent* —
  - Uses only one training sample or mini-batch to evaluate whether new parameter is more optimal than previous one
  - Fast and more likely to escape local minima due to randomness, but represents an approximation

- Algorithm:
  - Set  $\eta > 0$
  - Randomly initialize  $\beta_{(t=0)}$
  - Shuffle training data and initialize  $i \leftarrow 1$
  - $\beta_{(t+1)} \leftarrow \beta_{(t)} - \eta \nabla_\beta LO$  for observation  $i \mid \beta=\beta_{(t)}$
  - $t \leftarrow t + 1$
  - $i \leftarrow i + 1$
  - Repeat 4 to 6 until  $i = n + 1$
  - Repeat 2 to 6 until  $\nabla_\beta LO = 0$
- Basis for SGD is given by *Robbins-Monro algorithm*:
  - Algorithm:
    - Choose learning rates  $\eta_1, \eta_2, \dots$ , typically decreasing over time
    - Randomly initialize  $\beta_{(t=0)}$
    - $\beta_{(t+1)} \leftarrow \beta_{(t)} - \eta \nabla_\beta LO|_{\beta=\beta_{(t)}}$  where  $LO$  is noisy
  - For convergence:
    - $\sum_{t=1}^\infty \eta_t = \infty$  to ensure sufficient exploration
    - $\sum_{t=1}^\infty \eta_t^2 \leq \infty$  to avoid overly large updates
    - Then,  $\lim_{t \rightarrow \infty} P(|y_t - y| > \epsilon) = 0$  for any  $\epsilon > 0$

*Hyperparameters* —

- Learning rate  $\eta$ : Determines step size, if too small algorithm is slow to converge, if too large algorithm may diverge
- Batch size  $b$ : Number of samples from training set used to evaluate optimality of  $\beta$  at each step
- Epoch: Number of times model works through entire training set. Every epoch,  $\beta$  is updated  $n/b$  times

*Modifications* —

- Data should be standardized resp. scaled, otherwise the gradient of the largest predictor dominates the gradient of the loss function, leading to uneven updating of  $\beta$  and slow convergence
- A momentum term can be added to the updating function to ensure smooth updating of  $\beta$ :  $\beta_{(t+1)} \leftarrow \beta_{(t)} - \eta \nabla_\beta LO|_{\beta=\beta_{(t)}} + \alpha(\beta_{(t)} - \beta_{(t-1)})$
- For stochastic gradient descent, a smoothing step can be added because desired solution:  $\hat{\beta}_{(t+1)} \leftarrow \frac{1}{L+1} \sum_{j=t-L}^t \beta_{(t)}$

### 8 Model Evaluation

#### Estimator Evaluation Criteria

- Consistency:  $\hat{\theta} \rightarrow \theta$  as  $n \rightarrow \infty$
- Bias:  $\mathbb{E}(\hat{\theta}) - \theta$ 
  - Unbiased:  $\mathbb{E}(\hat{\theta}) = \theta$
  - Asymptotically unbiased:  $\mathbb{E}[(\hat{\theta} - \theta)^2] = 0$  as  $n \rightarrow \infty$
  - Asymptotically efficient:  $\mathbb{E}[(\hat{\theta} - \theta)^2] = I$  as  $n \rightarrow \infty$  where  $I$  is Fisher information (cf. Rao-Cramer bound)

#### Bias Variance Tradeoff

- Mean squared error  $\mathbb{E}[(\hat{f}(X) - y)^2]$  can be decomposed into:  $(\mathbb{E}[\hat{f}(X)] - f(X))^2 + \mathbb{V}(\hat{f}(X)) + \mathbb{E}[\epsilon^2] =$

bias<sup>2</sup> + variance + irreducible error

**Proof:**

- $y = f(X) + \epsilon$
- $\mathbb{E}[(\hat{f}(X) - y)^2] = \mathbb{E}[(\hat{f}(X) - \mathbb{E}[\hat{f}(X)] + \mathbb{E}[\hat{f}(X)] - f(X) + \epsilon)^2] = \mathbb{E}[(\hat{f}(X) - \mathbb{E}[\hat{f}(X)])^2] + \mathbb{E}[(\mathbb{E}[\hat{f}(X)] - f(X))^2] + \mathbb{E}[\epsilon^2] - 2\mathbb{E}[(\hat{f}(X) - \mathbb{E}[\hat{f}(X)])(\mathbb{E}[\hat{f}(X)] - f(X) + \epsilon)]$
- Fourth term equals 0:
  - $2\mathbb{E}[(\hat{f}(X) - \mathbb{E}[\hat{f}(X)])(\mathbb{E}[\hat{f}(X)] - f(X) + \epsilon)] = 2(\mathbb{E}[\hat{f}(X)] - f(X) + \epsilon)\mathbb{E}[(\hat{f}(X) - \mathbb{E}[\hat{f}(X)])]$  because  $(\mathbb{E}[\hat{f}(X)] - f(X) + \epsilon)$  is deterministic
  - In last equation, second term equals 0, so whole equation is 0
- Then, we are left with:
  - variance + bias<sup>2</sup> + irreducible error
- Bias: Error generated by the fact that we approximate a complex relationship via a simpler model (small function class) with a certain presupposed parametric form
- Variance: Error generated by the fact that we estimate the model parameters with a noisy training sample (small sample), rather than the population
- Irreducible error: Error generated by measurement error and the fact that we estimate  $y$  as a function of  $X$ , when it is a function of many other factors
- Bias variance tradeoff: Bias and variance cannot be reduced simultaneously
  - High variance associated with overfitting: Model corresponds too closely to particular training set resp. performs poorly on unseen data, but well on training set
  - High bias associated with underfitting: Model fails to capture underlying relationships resp. performs poorly on both training set and unseen data

#### Approximating Generalisation Loss via Empirical Loss

*Via resampling methods* —

Cross-validation:

- Partition data  $Z$  into  $K$  equally sized disjoint subsets:  $Z = Z_1 \cup Z_2 \cup \dots \cup Z_K$
- Produce estimator  $\hat{f}^{-v}$  from  $Z \setminus Z_v$  for  $v \leq K$
- Empirical loss given by:  $\hat{\mathcal{R}}^{cv} = \frac{1}{n} \sum_{i \leq n} LO(y_i - \hat{f}^{-k(i)}(x_i))$  where  $k(i)$  maps  $i$  to partition  $Z_{k(i)}$  where  $(x_i, y_i)$  belongs

Bootstrapping:

- Draw  $B$  samples with replacement of size  $n$  from data  $\mathcal{Z}$ :  $\mathcal{Z}^{*b}$
- Compute estimate  $S(\mathcal{Z}^{*b})$  for each bootstrap sample
- For each estimate, we can give a mean and variance:
  - $\bar{S} = \frac{1}{B} \sum_b S(\mathcal{Z}^{*b})$
  - $\sigma^2(S) = \frac{1}{B-1} \sum_b (S(\mathcal{Z}^{*b}) - \bar{S})^2$
- Empirical loss given by:  $\hat{\mathcal{R}}^{bs} = \frac{1}{n} \sum_{i \leq n} \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} LO(y_i - \hat{f}^{*b}(x_i))$  where  $C^{-i}$  contains all bootstrap indices  $b$

so that  $\mathcal{Z}^{*b}$  does not contain  $(x_i, y_i)$

- Empirical loss of bootstrap uses training data to estimate  $\hat{\mathcal{R}}$ , i.e. it is generally too optimistic. We can correct this:
  - Probability that  $(x_i, y_i)$  is not in sample  $\mathcal{Z}^{*b}$  of size  $n$  is given by  $(1 - \frac{1}{n})^n = \frac{1}{e}$  as  $n \rightarrow \infty \approx \frac{1}{3}$
  - Probability that  $(x_i, y_i)$  is in sample  $\mathcal{Z}^{*b}$  of size  $n$  is given by  $1 - \frac{1}{e}$  as  $n \rightarrow \infty \approx \frac{2}{3}$
  - We then define:  $\hat{\mathcal{R}}^{(0.632)} = 0.368\hat{\mathcal{R}} + 0.632\hat{\mathcal{R}}^{bs}$

### 9 Estimating Common Distributions

#### Gaussian

*Frequentism (MLE)* —

- Likelihood (excl. constants):  $L = (\frac{1}{\sigma})^n \prod_{i=2}^n \exp(-\frac{1}{2\sigma^2}(x^{(i)} - \mu)^2)$
- Log-likelihood:  $LL = -n \log(\sigma) - \sum_{i=1}^n (\frac{1}{2\sigma^2}(x^{(i)} - \mu)^2)$
- $\mu_{MLE}$  is sample mean:  $\frac{1}{n} \sum_{i=1}^n x^{(i)}$ :
  - Derivative of log-likelihood wrt  $\mu$ :

$$\begin{aligned} \nabla_\mu LL &= \nabla_\mu - \sum_{i=1}^n (\frac{x^{(i)2} - 2x^{(i)}\mu + \mu^2}{2\sigma^2}) = \\ \nabla_\mu - \sum_{i=1}^n (-\frac{x^{(i)}\mu}{\sigma^2} + \frac{\mu^2}{2\sigma^2}) &= -\sum_{i=1}^n (-\frac{x^{(i)}}{\sigma^2} + \frac{2\mu}{2\sigma^2}) = \sum_{i=1}^n (\frac{x^{(i)} - \mu}{\sigma^2}) = \sum_{i=1}^n x^{(i)} - n\mu = 0 \end{aligned}$$

- $\sigma^2_{MLE}$  is sample variance:  $\frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)^2$ :
  - Derivative of log-likelihood wrt  $\sigma$ :  $\nabla_\sigma LL =$ 

$$\begin{aligned} -n \nabla_\sigma \log(\sigma) - \nabla_\sigma (\sum_{i=1}^n (\frac{(x^{(i)} - \mu)^2}{2\sigma^2})) &= \\ -\frac{n}{\sigma} - \nabla_\sigma (\sum_{i=1}^n \frac{1}{2} \sigma^{-2} (x^{(i)} - \mu)^2) &= \\ -\frac{n}{\sigma} - (\sum_{i=1}^n -1 \sigma^{-3} (x^{(i)} - \mu)^2) &= \\ -n + \sum_{i=1}^n (\frac{(x^{(i)} - \mu)^2}{\sigma^2}) &= 0 \end{aligned}$$

*Bayesianism* —

- Assume  $\Sigma$  is known and  $\mu \sim \mathcal{N}(\mu_0, \Sigma_0)$  is the outcome of a random variable
- $p(\mu|X, \mu_0, \Sigma_0) \propto p(X|\mu, \Sigma)p(\mu|\mu_0, \Sigma_0)$
- $p(X|\mu, \Sigma) = \frac{1}{2\pi^{mn/2}} \frac{1}{|\Sigma|^{n/2}} \exp(\frac{1}{2} \sum_{i=1}^n (x^{(i)} - \mu)^\top \Sigma^{-1} (x^{(i)} - \mu))$
- $p(\mu|\mu_0, \Sigma_0) = \frac{1}{2\pi^{m/2}} \frac{1}{|\Sigma_0|^{m/2}} \exp(\frac{1}{2} \sum_{i=1}^n (\mu - \mu_0)^\top \Sigma_0^{-1} (\mu - \mu_0))$
- $p(\mu|X, \mu_0, \Sigma_0) \propto \exp(-\frac{1}{2}(\mu^\top \Sigma_0^{-1} \mu + n\mu^\top \Sigma^{-1} \mu - 2\mu_0^\top \Sigma_0^{-1} \mu - 2n\bar{x}^\top \Sigma^{-1} \mu))$  after combining exponents of the prior and likelihood, expanding, absorbing terms unrelated to  $\mu$  into a constant, and replacing  $\sum_{i=1}^n x^{(i)\top}$  by  $n\bar{x}^\top$
- We now apply a symmetric matrix property  $x^\top Ax + 2x^\top b = (x + A^{-1}b)^\top A(x + A^{-1}b) - b^\top A^{-1}b$ , with  $\mu = x$ ,  $-(\Sigma_0^{-1} + n\Sigma^{-1})^{-1} = A^{-1}$  and  $(\Sigma^{-1}n\bar{x} + \Sigma_0^{-1}\mu_0) = b$
- Through this, we get  $p(\mu|X, \mu_0, \Sigma_0) \propto \exp(\frac{1}{2}(\mu(\Sigma_0^{-1} + n\Sigma^{-1})^{-1}(\Sigma^{-1}n\bar{x} + \Sigma_0^{-1}\mu_0))^\top (\Sigma_0^{-1} +$



- $n\Sigma^{-1})(\mu - (\Sigma_0^{-1} + n\Sigma^{-1})^{-1}(\Sigma^{-1}n\bar{x} + \Sigma_0^{-1}\mu_0))) = \exp(\frac{1}{2}(\mu - \mu_n)^\top \Sigma_n^{-1}(\mu - \mu_n))$
- Thus,  $p(\mu|X, \mu_0, \Sigma_0) \sim \mathcal{N}(\mu_n, \Sigma_n)$  with
  - $\mu_n = (\Sigma_0^{-1} + n\Sigma^{-1})^{-1}(\Sigma^{-1}n\bar{x} + \Sigma_0^{-1}\mu_0) =$  (if  $\Sigma$  equals 1)  $\frac{n\bar{x}\Sigma_0 + \mu_0}{n\Sigma_0 + 1}$
  - $\Sigma_n = (\Sigma_0^{-1} + n\Sigma^{-1})^{-1} =$  (if  $\Sigma$  equals 1)  $\frac{\Sigma_0}{n\Sigma_0 + 1}$
- For Bayesian parameter  $\mu_n$ :
  - $\mu_n$  is a compromise between MLE and prior, approximating prior for small  $n$  and MLE for large  $n$
  - If prior variance is small (i.e. if we are certain of our prior), prior mean weighs more strongly
- For Bayesian parameter  $\Sigma_n$ :
  - $\Sigma_n$  approximates prior for small  $n$  and MLE for large  $n$
  - If prior variance is small (i.e. if we are certain of our prior), posterior variance is also small

#### Binomial

#### TBA

#### Poisson

#### TBA

### 10 Linear Regression

#### Description

*Task* — Regression

*Description* —

- Supervised
- Parametric

#### Formulation

- $y^{(i)} = \beta \cdot x^{(i)}$  resp.  $y = X\beta$  where  $X$  contains  $n$  rows, each of which represents an instance, and  $m$  columns, each of which represents a feature
- This can be considered as a projection of  $y$  to the columnspace of  $X$

#### Optimization

*Parameters* — Find parameters  $\beta$

*Objective function* — Ordinary least squares estimator (OLSE):

- Minimize mean squared error:
 
$$LO = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \beta \cdot x^{(i)})^2$$
 resp.  $LO = (y - X\beta)^\top (y - X\beta)$

MLE:

- Yields same result as OLSE
- Orthogonality principle:
  - Yields same result as OLSE
  - $\hat{y} = X\beta$  is a projection of  $y$  to the columnspace of  $X$
- Then, by the orthogonality principle,  $X \cdot (\hat{y} - y) = X \cdot (X\beta - y) = 0$

- $\Rightarrow \beta = (X^\top X)^{-1} X^\top y$

*Optimization* —

- $\nabla_\beta LO = \frac{1}{2} \nabla_\beta ((y - X\beta)^\top (y - X\beta)) = \frac{1}{2} \nabla_\beta (\beta^\top X^\top X \beta - 2y^\top X \beta) = X^\top X \beta - X^\top y = X^\top (X\beta - y) = 0$
- $\Rightarrow \beta = (X^\top X)^{-1} X^\top y$

*Hypothesis Testing of Found Parameters* —

- Let  $y|X \sim \mathcal{N}(y, \sigma^2 I) = \mathcal{N}(X\beta, \sigma^2 I)$
- Let  $\hat{\beta} = (X^\top X)^{-1} X^\top y = X^+ y$  be the OLSE where  $X^+$  is a scalar

- Then,  $\hat{\beta} \sim \mathcal{N}(X^+ X \beta, X^+ \tau \sigma^2 X^+) = \mathcal{N}(\beta, (X^\top X)^{-1} \sigma^2)$ 

**Proof:**

  - $\mathcal{N}(X^+ X \beta, X^+ \tau \sigma^2 X^+) = \mathcal{N}(I \beta, \sigma^2 X^+ (X^\top X)^{-1} X^\top) = \mathcal{N}(I \beta, \sigma^2 X^+ (X^\top X)^{-1} X^\top) = \mathcal{N}(\beta, \sigma^2 X^+ X (X^\top X)^{-1} \tau) = \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1})$  since  $(X^\top X)$  is symmetric

- We can estimate  $\sigma^2$  unbiasedly as:
 
$$\hat{\sigma}^2 = \frac{1}{n-m} \sum_{i \leq n} (X\hat{\beta} - y)^2$$
- Then, confidence interval for  $\hat{\beta}_j$  given by:  $\hat{\beta}_j \pm z_{\alpha/2} se(\hat{\beta}_j)$  where
  - $z_{\alpha/2} = \Phi^{-1}(\alpha/2)$  is Gaussian CDF
  - $se(\hat{\beta}_j)$  is the  $j^{th}$  diagonal element of the covariance matrix  $\sigma^2 (X^\top X)^{-1}$
- We can perform a hypothesis test on  $\hat{\beta}$  with the *Wald test*:
  - $H_0: \beta = \beta_0$  (typically 0)
  - $H_1: \beta \neq \beta_0$
  - Wald statistic:  $W = \frac{\hat{\beta} - \beta_0}{\hat{se}}$
  - If p-value associated with  $W$  is smaller than  $\alpha$  resp. if  $|W|$  is greater than or equal to the critical value  $z_{\alpha/2}$ , we reject  $H_0$

*Evaluation* —

- OLSE is unbiased if noise  $\epsilon$  has zero mean:
  - Given  $y = X\beta + \epsilon$ , we can substitute  $\hat{\beta} = (X^\top X)^{-1} X^\top (X\beta + \epsilon) = \beta + (X^\top X)^{-1} X^\top \epsilon$
  - Taking the expected value on both sides, we have:
 
$$\mathbb{E}(\hat{\beta}) = \beta + (X^\top X)^{-1} X^\top \mathbb{E}(\epsilon)$$
  - Then,  $\mathbb{E}(\hat{\beta}) = \beta$  if the noise has zero mean
- Gauss Markov theorem*: OLSE is best (lowest variance, lowest MSE) unbiased estimator, if assumptions ( $X$  is full rank and there is no multicollinearity, heteroskedasticity, and exogeneity) are met

- Proof:**
- Let  $A^\top y = (X^\top X)^{-1} X^\top y$  be the OLSE
  - Let  $C^\top y$  be another unbiased estimator
  - $\mathbb{V}(A^\top y) = A^\top \mathbb{V}(y) A$  since  $A$  is constant
  - We can further develop to:
 
$$A^\top \sigma^2 I_m A = \sigma^2 A^\top A$$
 since variance is given by error term
    - Similarly,  $\mathbb{V}(C^\top y) = \sigma^2 C^\top C$
    - For the OLSE, we can plug in  $(X^\top X)^{-1} X^\top$  for  $A$  which yields:
 
$$\mathbb{V}(A^\top y) = \sigma^2 (X^\top X)^{-1} X^\top X (X^\top X)^{-1} = \sigma^2 (X^\top X)^{-1}$$
    - Then, we have shown that  $\mathbb{V}(A^\top y) \leq \mathbb{V}(C^\top y)$

- Nonetheless, there may be biased estimators that generate a lower variance and MSE

*Characteristics* —

- Convex with psd Hessian
- Has global minimum
- Has unique or infinitely many solutions
- Can be solved analytically, if  $(X^\top X)$  is

invertible

#### 11 Linear Minimum Mean Squared Error Estimation (LMMSE)

#### Description

- Minimizes mean squared error of two random variables, leveraging information about their mean and covariance
- Linear regression with large samples approximates LMMSE

#### Formulation

- $y$  is observed
- $x$  is a row vector and quantity of interest
- We estimate  $x$  as  $\hat{x} = h^\top Y = \sum_i h_i y_i$  where  $X$  contains  $m$  rows, each of which represents the  $n$ -sized vector for a random variable
- This can be considered as a projection of  $x$  to the rowspace of  $Y$

#### Optimization

*Parameters* — Find parameters

*Objective function* —

- Minimize expected squared error:
 
$$LO = \mathbb{E}[|\hat{x} - x|]$$

*Optimization* —

- By the orthogonality principle,  $\mathbb{E}[(\hat{x} - x) \cdot y_i] = \mathbb{E}[(\sum_{l=1}^n h_l y_l - x) \cdot y_i] = 0$  for  $i = 1, \dots, n$
- Then,  $\sum_{l=1}^n \mathbb{E}[y_l \cdot y_i] h_l = \mathbb{E}[x \cdot y_i]$  for  $i = 1, \dots, n$  which in matrix notation corresponds to
 
$$\begin{bmatrix} \mathbb{E}[y_1 \cdot y_1] & \dots & \mathbb{E}[y_1 \cdot y_n] \\ \mathbb{E}[y_n \cdot y_1] & \dots & \mathbb{E}[y_n \cdot y_n] \\ \mathbb{E}[x \cdot y_1] \\ \vdots \\ \mathbb{E}[x \cdot y_n] \end{bmatrix} \text{ resp. concisely } h^\top \mathbb{E}[Y Y^\top] = \mathbb{E}[x Y^\top]$$

### 12 Bayesian Linear Regression

#### Description

*Task* — Regression

*Description* —

- Supervised
- Parametric

#### Formulation

- $y^{(i)} = \beta \cdot x^{(i)}$  resp.  $y = X\beta$
- $\beta \sim \mathcal{N}(0, T^2 I_m)$
- $p(\beta) \propto -\frac{1}{2T^2} \beta^\top \beta$

#### Optimization

*Parameters* — Find distribution of

parameters  $\beta$

*Optimization* —

- Likelihood:
  - Conditional on  $\beta$ ,  $y \propto \mathcal{N}(X\beta, \sigma^2 I_m)$
  - $p(y|X, \beta) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp(-\frac{1}{2\sigma^2} (y - X\beta)^\top (y - X\beta))$
- Posterior:
 
$$p(\beta|X, y) \propto p(X, \beta) \times p(\beta) \propto \exp(-\frac{1}{2\sigma^2} (y - X\beta)^\top (y - X\beta)) \times \exp(-\frac{1}{2T^2} \beta^\top \beta) = \exp(-\frac{1}{2} (\frac{1}{\sigma^2} y^\top y - 2\beta^\top X^\top y + \beta^\top X^\top X \beta) + \frac{1}{2T^2} \beta^\top \beta) \propto \exp(-\frac{1}{2} (\beta^\top (\frac{1}{\sigma^2} X^\top X + \frac{1}{2T^2} I_m) \beta - \frac{2}{\sigma^2} \beta^\top X^\top y))$$
- We now apply a symmetric matrix property  $x^\top A x + 2x^\top b = (x + A^{-1} b)^\top A (x + A^{-1} b) - b^\top A^{-1} b$ , with

$$\beta = x, (\frac{1}{\sigma^2} X^\top X + \frac{1}{2T^2} I_m) = A \text{ and}$$

$$(\frac{1}{\sigma^2} X^\top y) = b$$

- Through this, we get  $p(\beta|X, y) \propto \exp(\frac{1}{2}(\beta + (\frac{1}{\sigma^2} X^\top X + \frac{1}{T^2} I_m)^{-1} (\frac{1}{\sigma^2} X^\top y))^\top (\frac{1}{\sigma^2} X^\top X + \frac{1}{T^2} I_m) (\beta + (\frac{1}{\sigma^2} X^\top X + \frac{1}{T^2} I_m)^{-1} (\frac{1}{\sigma^2} X^\top y)))$
- Thus,  $p(\beta|X, y) \sim \mathcal{N}(\mu, \Sigma)$  with
  - $\mu = \Sigma \times \frac{1}{\sigma^2} X^\top y$
  - $\Sigma = (\frac{1}{\sigma^2} X^\top X + \frac{1}{T^2} I_m)^{-1}$
- Posterior mean corresponds to parameter  $\beta$  found by ridge regression, if  $\lambda = \frac{\sigma^2}{T^2}$
- If we set an infinitely broad prior  $T^2$  then the Bayesian estimate converges to the MLE estimate – if we have  $n = 0$  training instances, the Bayesian estimate reverts to the prior

*Characteristics* —

- Convex with psd Hessian
- Has global minimum
- Can be solved analytically

### 13 Ridge ( $\ell_2$ ) Regression

#### Description

*Task* — Regression

*Description* —

- Supervised
- Parametric

#### Formulation

- $y^{(i)} = \beta \cdot x^{(i)}$  resp.  $y = X\beta$

#### Optimization

*Parameters* — Find parameters  $\beta$  subject to

$$\|\beta\|^2 \leq t \text{ resp. } \|\beta\|^2 - t \leq 0$$

*Objective function* —

- Minimize mean squared error subject to constraint
- Lagrangian formulation:
 
$$LO = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \beta \cdot x^{(i)})^2 + \lambda (\|\beta\|^2 - t)$$
 resp.  $LO = (y - X\beta)^\top (y - X\beta) + \lambda (\|\beta\|^2 - t)$
- Still a OLSE problem, since we can rewrite the objective to minimize  $(X\beta - y)$  as the objective to minimize  $\|(X' \beta - y')\|^2$  with  $X' = \begin{bmatrix} X \\ \lambda I \end{bmatrix}$  and  $y' = \begin{bmatrix} y \\ 0 \end{bmatrix}$

*Optimization* —

- $\nabla_\beta LO = 0$

$$\Rightarrow \beta = (X^\top X + \lambda I)^{-1} X^\top y$$

*Effect* —

- Shrinks certain elements of  $\beta$  to near 0
 

**Proof:**

  - Gradient at optimality given by  $\frac{\partial (y - X\beta)^\top (y - X\beta)}{\partial \beta} + 2\lambda \beta = 0$
  - Then,  $\beta^* = -\frac{1}{2\lambda} \frac{\partial (y - X\beta)^\top (y - X\beta)}{\partial \beta}$
  - This means that each parameter is shrunk by a factor determined by size of  $\lambda$  - the larger  $\lambda$ , the more the parameters are shrunk
  - Larger parameters experience a larger shrinkage

*Characteristics* —

- Strictly with pd Hessian, since Lagrangian term is strictly convex and the sum of a strictly convex function with a convex function is strictly convex

- Has global minimum
- Has unique solution, as  $(X^\top X + \lambda I)$  has linearly independent columns
- Can be solved analytically, as  $(X^\top X + \lambda I)$  is always invertible

### 14 Lasso ( $\ell_1$ ) Regression

#### Description

*Task* — Regression

*Description* —

- Supervised
- Parametric

#### Formulation

- $y^{(i)} = \beta \cdot x^{(i)}$  resp.  $y = X\beta$

#### Optimization

*Parameters* — Find parameters  $\beta$  subject to  $|\beta| \leq t$  resp.  $|\beta| - t \leq 0$

*Objective function* —

- Minimize mean squared error subject to constraint
- Lagrangian formulation:
 
$$LO = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \beta \cdot x^{(i)})^2 + \lambda (|\beta| - t)$$
 resp.  $LO = (y - X\beta)^\top (y - X\beta) + \lambda (|\beta| - t)$

*Effect* —

- Shrinks certain elements of  $\beta$  to 0
 

**Proof:**

  - Gradient at optimality given by  $\frac{\partial (y - X\beta)^\top (y - X\beta)}{\partial \beta} + \frac{\partial \lambda |\beta|}{\partial \beta} = 0$
  - $\frac{\partial \lambda |\beta|}{\partial \beta}$  non-differentiable because there is a sharp edge at  $\beta = 0$ , but we can work with subgradients for  $\beta \neq 0$
  - If we have  $-\lambda < \frac{\partial (y - X\beta)^\top (y - X\beta)}{\partial \beta} < \lambda$  the optimum is given by  $\beta = 0$
  - This means that some parameters are set to 0 - the larger  $\lambda$ , the more parameters are set to 0
  - Small parameter values (i.e. unimportant features) are more likely to be set to 0
  - For parameters that are not set to 0, LASSO regression has a similar effect as ridge regression and shrinks these parameters towards 0

*Characteristics* —

- Convex, but not strictly convex
- Has global minimum
- Has unique or infinitely many solutions
- Cannot be solved analytically, since  $|\beta|$  is not differentiable at  $\beta_i = 0$

### 15 Polynomial Regression

#### Description

*Task* — Regression

*Description* —

- Supervised
- Parametric

#### Formulation

- $y^{(i)} = \beta \cdot \phi(x^{(i)})$  resp.  $y = \Phi \beta$  where  $\Phi$  is the transformed design matrix with rows  $\phi(x^{(i)})^\top$

<b>Optimization</b>
<i>Parameters</i> — Find parameters $\beta$
<i>Objective function</i> —
<ul style="list-style-type: none"> <li>Ordinary least squares estimator</li> <li>Minimize mean squared error</li> </ul>
<i>Optimization</i> —
<ul style="list-style-type: none"> <li><math>\nabla_{\beta} LO = 0</math></li> <li><math>\Rightarrow \beta = (\Phi^{\top} \Phi)^{-1} \Phi^{\top} y</math></li> </ul>
<i>Characteristics</i> —
<ul style="list-style-type: none"> <li>Convex with psd Hessian</li> <li>Has global minimum</li> <li>Has unique or infinitely many solutions</li> <li>Can be solved analytically, if <math>(\Phi^{\top} \Phi)</math> is invertible</li> </ul>
<b>16 Kernel Methods</b>
<b>Background on Kernel Methods</b>
<i>Description</i> —
<ul style="list-style-type: none"> <li>Mechanism for tractably resp. implicitly mapping data into higher-dimensional feature space so that linear models can be used in this feature space</li> <li>To do so, we can employ the <i>kernel trick</i> and the <i>representer theorem</i></li> <li>The requirements are that the kernel function fulfills <i>Mercer's theorem</i>, i.e. the kernel is a Mercer kernel</li> </ul>
<i>Kernel trick</i> —
<ul style="list-style-type: none"> <li>Allows to operate in higher-dimensional feature space, without explicitly calculating this transformation, but instead implicitly computing the inner product in this feature space via a kernel function</li> </ul>
<ul style="list-style-type: none"> <li>Given two inputs <math>x^{(i)}, x^{(j)}</math> and a feature map <math>\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^k</math> we can define an inner product on <math>\mathbb{R}^k</math> via the kernel function: <math>k(x^{(i)}, x^{(j)}) = \varphi(x^{(i)}) \cdot \varphi(x^{(j)})</math></li> </ul>
<ul style="list-style-type: none"> <li>If a prediction function is described solely in terms of inner products in the input space, it can be lifted into the feature space by replacing the inner product with the kernel function</li> <li>Kernel trick cannot be used in conjunction with feature selection resp. sparsity inducing regularize (e.g. <math>\ell_1</math>), as this does not satisfy the representer theorem</li> </ul>
<i>Representer theorem</i> —
<ul style="list-style-type: none"> <li>Allows to avoid directly seeking the <math>k</math> parameters, but only the <math>n</math> parameters that characterize <math>\alpha</math></li> <li>Allows to avoid calculating <math>\varphi(z)</math> when evaluating novel instance, but only sum over weighted set of <math>n</math> kernel function outputs</li> </ul>
<i>Mercer's theorem</i> —
<ul style="list-style-type: none"> <li>Kernel function is psd and symmetric iff <math>k(x^{(i)}, x^{(j)}) = \varphi(x^{(i)}) \cdot \varphi(x^{(j)})</math> <ul style="list-style-type: none"> <li>Psd: <math>x^{\top} K x \geq 0</math> where <math>K</math> is the kernel matrix</li> <li>Symmetric: <math>k(x^{(i)}, x^{(j)}) = k(x^{(j)}, x^{(i)})</math></li> </ul> </li> <li>Kernel that satisfies Mercer's theorem is a Mercer kernel, i.e. we can prove a kernel is a Mercer kernel either if it is psd and symmetric or by finding a feature map such that the kernel function corresponds to an inner product</li> </ul>

<b>Formulation</b>
<ul style="list-style-type: none"> <li>Feature map <math>\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^k</math></li> <li>Linear prediction function: <math>\beta \cdot \varphi(x^{(i)})</math></li> <li>Regularized loss function: <math>LO = \sum_{i=1}^n LO(y^{(i)}, \beta \cdot \varphi(x^{(i)}) + \Omega(\beta))</math></li> <li>Iff <math>\Omega(\beta)</math> is a non-decreasing function, then the parameters <math>\beta</math> that minimize the loss function can be rewritten as: <math>\beta = \sum_{i=1}^n \alpha^{(i)} \varphi(x^{(i)})</math></li> <li>Outcome of novel instance can be predicted as: <math>\beta \cdot \varphi(z) = \sum_{i=1}^n \alpha^{(i)} \varphi(x^{(i)}) \cdot \varphi(z) = \sum_{i=1}^n \alpha^{(i)} k(x^{(i)}, z)</math></li> <li>Act of prediction becomes act of measuring similarity to training instances in feature map space</li> </ul>

<b>Kernel Types</b>
<i>Polynomial kernel</i> —
<ul style="list-style-type: none"> <li><math>\varphi(x) = [x^{\alpha}]_{\alpha \in \mathbb{N}^m}</math> where <math>\alpha = (\alpha_1, \dots, \alpha_m)</math> is the multi-index representing the power and <math>x^{\alpha} = x_1^{\alpha_1} \times \dots \times x_m^{\alpha_m}</math> is the monomial term corresponding to the multi-index <math>\alpha</math></li> <li>E.g. if degree = 2, then <math>k(x^{(i)}, x^{(j)}) = 1 + 2x_1^{(i)} x_1^{(j)} + 2x_2^{(i)} x_2^{(j)} + (x_1^{(i)} x_1^{(j)})^2 + (x_2^{(i)} x_2^{(j)})^2 + 2x_1^{(i)} x_1^{(j)} x_2^{(i)} x_2^{(j)}</math></li> <li>Inner product diverges to infinity</li> <li>To address this, we often use RBF kernel instead</li> </ul>
<i>RBF kernel</i> —
<ul style="list-style-type: none"> <li>Gives access to infinite feature space</li> <li><math>\varphi(x) = \exp(-\frac{1}{2} \ x\ ^2) [\frac{x^{\alpha}}{\sqrt{\alpha!}}]_{\alpha \in \mathbb{N}^m}</math></li> <li><math>k(x^{(i)}, x^{(j)}) = \sigma^2 \exp(-\frac{\ x^{(i)} - x^{(j)}\ ^2}{2l^2})</math> <ul style="list-style-type: none"> <li>Proof: <math>-\exp(-\frac{1}{2} \ x^{(i)}\ ^2) \exp(-\frac{1}{2} \ x^{(j)}\ ^2) \sum_{\alpha} [\frac{x^{(i)\alpha} x^{(j)\alpha}}{\alpha!}]</math></li> <li>Given multinomial series expansion, <math>\sum_{\alpha} [\frac{x^{(i)\alpha} x^{(j)\alpha}}{\alpha!}] = \exp(x^{(i)\top} x^{(j)})</math></li> <li><math>-\exp(-\frac{1}{2} \ x^{(i)}\ ^2 - \frac{1}{2} \ x^{(j)}\ ^2 + x^{(i)\top} x^{(j)}) = \exp(-\frac{\ x^{(i)} - x^{(j)}\ ^2}{2\sigma^2})</math></li> </ul> </li> <li>Length scale parameter <math>l</math> controls how quickly the similarity decays with distance</li> <li>Variance parameter <math>\sigma</math> controls the vertical scale of the function</li> </ul>

<i>Kernel compositions</i> —
<ul style="list-style-type: none"> <li>New valid kernels can be composed via: <ul style="list-style-type: none"> <li>Addition: <math>k_1 + k_2</math></li> <li>Multiplication: <math>k_1 \times k_2</math></li> <li>Scaling: <math>c \times k_1</math> for <math>c &gt; 0</math></li> <li>Composition: <math>f(k_1)</math> where <math>f</math> is a polynomial with positive coefficients or the exponential function</li> </ul> </li> </ul>

<b>17 Polynomial Kernel Regression</b>
<b>Description</b>
<i>Task</i> — Regression
<i>Description</i> —
<ul style="list-style-type: none"> <li>Supervised</li> <li>Parametric</li> </ul>
<b>Formulation</b>
<ul style="list-style-type: none"> <li><math>y = \beta \cdot \varphi(x^{(i)})</math></li> </ul>

<b>Optimization</b>
<i>Parameters</i> — Find parameters $\beta$
<i>Objective function</i> —
<ul style="list-style-type: none"> <li>Ordinary least squares estimator (OLSE)</li> <li>Minimize mean squared error: <math>LO = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \beta \cdot \varphi(x^{(i)}))^2</math></li> </ul>
<i>Optimization</i> —
<ul style="list-style-type: none"> <li>Primal solution: <ul style="list-style-type: none"> <li>Parameters can be estimated as: <math>\beta = (\Phi^{\top} \Phi)^{-1} \Phi^{\top} y</math></li> <li>Prediction for novel instance: <math>\beta \cdot \varphi(z) = (\Phi^{\top} \Phi)^{-1} \Phi^{\top} y \cdot \varphi(z) = y^{\top} \Phi (\Phi^{\top} \Phi)^{-1} \varphi(z)</math></li> </ul> </li> <li>Let us define <math>K = \Phi \Phi^{\top}</math> as the kernel matrix of the training data with <math>K_{ij} = \varphi(x^{(i)}) \cdot \varphi(x^{(j)})</math></li> </ul>

<ul style="list-style-type: none"> <li>Dual solution <math>\alpha</math> if we have no regularization, i.e. <math>\lambda = 0</math>: <ul style="list-style-type: none"> <li>Parameters can be estimated as: <math>\beta = \Phi^{\top} K^{-1} y</math> <ul style="list-style-type: none"> <li>Proof: <ul style="list-style-type: none"> <li><math>* (\Phi^{\top} \Phi + \lambda I) \beta = \Phi^{\top} y</math></li> <li><math>* \Rightarrow \Phi^{\top} \Phi \beta + \lambda I \beta = \Phi^{\top} y</math></li> <li><math>* \Rightarrow I \beta = \Phi^{\top} \lambda^{-1} (y - \Phi \beta)</math></li> <li><math>* \text{ Since we know from the representer theorem that } \beta = \Phi^{\top} \alpha, \text{ we can say: } \alpha = \lambda^{-1} (y - \Phi \beta)</math></li> <li><math>* \text{ We can further develop this to: } \lambda \alpha = (y - \Phi \beta)</math></li> <li><math>* \text{ Replacing } \beta \text{ by } \Phi^{\top} \alpha \text{ yields: } \lambda \alpha = (y - \Phi \Phi^{\top} \alpha)</math></li> <li><math>* \Rightarrow \alpha = (\Phi \Phi^{\top} + \lambda I)^{-1} y = K^{-1} y</math></li> <li><math>* \text{ With this, we can calculate the parameters: } \beta = \Phi^{\top} \alpha = \Phi^{\top} (\Phi \Phi^{\top} + \lambda I)^{-1} y = \Phi^{\top} K^{-1} y</math></li> </ul> </li> </ul> </li> <li>Prediction for novel instance: <math>\beta \cdot \varphi(z) = y^{\top} (\Phi \Phi^{\top})^{-1} \Phi \varphi(z) = y^{\top} (\Phi \Phi^{\top})^{-1} k</math> where <math>k = \Phi \varphi(z) = [k(x^{(1)}, z), \dots, k(x^{(n)}, z)]^{\top} = [\varphi(x^{(1)}) \cdot \varphi(z), \dots, \varphi(x^{(n)}) \cdot \varphi(z)]^{\top}</math> is a kernel vector, consisting of kernel values between training instances and new instance</li> </ul> </li></ul>
--

<i>Algorithm</i> — Training:
<ol style="list-style-type: none"> <li>Compute kernel matrix given RBF kernel Time complexity: <math>\mathcal{O}(n^2 \times m)</math> for <math>n^2</math> kernel matrix values and <math>m</math> number of features in each instance vector</li> <li>Perform training by solving <math>\alpha = K^{-1} y</math> for <math>\alpha</math> Time complexity: <math>\mathcal{O}(n^3)</math></li> <li>Store <math>\alpha</math> Space complexity: <math>\mathcal{O}(n^2)</math></li> </ol>
Prediction:
<ol style="list-style-type: none"> <li>Compute kernel vector Time complexity: <math>\mathcal{O}(n \times m \times d)</math> for <math>d</math> new instances, given <math>n</math> instances in training data and <math>m</math> features in each instance vector</li> <li>Store <math>k</math> Space complexity: <math>\mathcal{O}(n \times d)</math> for <math>d</math> new instances, given <math>n</math> as length of kernel vector</li> <li>Predict response using stored kernel vector Time complexity: <math>\mathcal{O}(n \times d)</math> for <math>d</math> new instances, given <math>n</math> as length of <math>\alpha</math></li> </ol>

Value:
<ul style="list-style-type: none"> <li>Primal solution training is of time complexity <math>\mathcal{O}(k^3)</math> and prediction is of time complexity <math>\mathcal{O}(k)</math></li> <li>Dual solution speeds this up as seen above in the algorithm</li> </ul>

<i>Characteristics</i> —
<ul style="list-style-type: none"> <li>Convex with psd Hessian</li> <li>Has global minimum</li> <li>Has unique or infinitely many solutions</li> <li>Can be solved analytically</li> </ul>

<b>18 Gaussian Processes</b>
<b>Description</b>
<i>Task</i> — Models a distribution over functions.
<i>Description</i> —
<ul style="list-style-type: none"> <li>Supervised</li> <li>Non-parametric</li> </ul>
<b>Formulation</b>
<ul style="list-style-type: none"> <li><math>y^{(i)} = \beta \cdot x^{(i)} + \epsilon</math> resp. <math>y = X \beta + \epsilon</math></li> <li><math>\beta \sim \mathcal{N}(0, \Lambda^{-1})</math></li> <li><math>\epsilon \sim \mathcal{N}(0, \sigma I_m)</math></li> </ul>

<b>Optimization</b>
<i>Optimization</i> —
<ul style="list-style-type: none"> <li>If we compute the moment of the Gaussian: <ul style="list-style-type: none"> <li><math>\mathbb{E}[y] = X^{\top} \mathbb{E}(\beta) = X^{\top} 0 = 0</math></li> <li><math>\text{Cov}(y) = \mathbb{E}[(X^{\top} \beta + \epsilon)(X^{\top} \beta + \epsilon)^{\top}] = X \mathbb{E}(\beta \beta^{\top}) = X^{\top} + X \mathbb{E}(\beta) \mathbb{E}(\epsilon^{\top}) + \mathbb{E}(\epsilon) \mathbb{E}(\beta^{\top}) X^{\top} + \mathbb{E}(\epsilon \epsilon^{\top}) = 0</math> where <math>* \mathbb{V}(\beta) = \mathbb{E}(\beta \beta^{\top})</math> and <math>\mathbb{V}(\epsilon) = \mathbb{E}(\epsilon \epsilon^{\top})</math> because <math>\mathbb{V}(x) = \mathbb{E}[(x - \mathbb{E}(x))^2] = \mathbb{E}[x^2]</math> if <math>\mathbb{E}(x) = 0</math>, which is the case here due to the defined distributions</li> <li><math>* \mathbb{E}(\epsilon) = 0</math></li> </ul> </li> <li>Plugging in the variance for <math>\beta</math> and <math>\epsilon</math>, we have <math>\text{Cov}(y) = X \Lambda^{-1} X^{\top} + \sigma^2 I_m</math></li> <li>This can be written as a Kernel matrix <math>K</math>: <math display="block">\begin{bmatrix} K_{1,1} + \sigma^2 &amp; \dots &amp; \dots &amp; K_{1,n} \\ \dots &amp; K_{2,2} + \sigma^2 &amp; \dots &amp; \dots \\ \dots &amp; \dots &amp; \dots &amp; \dots \\ K_{n,1} &amp; \dots &amp; \dots &amp; K_{n,n} + \sigma^2 \end{bmatrix}</math> with <math>K_{ij} = x^{(i)\top} \Lambda^{-1} x^{(j)}</math></li> <li>In this kernel matrix, the kernel function can take any shape</li> </ul>

<ul style="list-style-type: none"> <li>On this basis, Gaussian process is defined as collection of random variables such that every finite subset of variables is jointly Gaussian: <math>f \sim \mathcal{GP}(\mu, K)</math></li> <li>A new instance follows the distribution <math>p(y_{n+1}) = \mathcal{N}(k^{\top} C_n^{-1} y, c - k^{\top} C_n^{-1} k)</math> where <ul style="list-style-type: none"> <li><math>k = k(x^{(1)}, x^{(n+1)}), \dots, k(x^{(n)}, x^{(n+1)})</math></li> </ul> <math>]^{\top} = [\varphi(x^{(1)}) \cdot \varphi(x^{(n+1)}), \dots, \varphi(x^{(n)}) \cdot \varphi(x^{(n+1)})]^{\top}</math> is the kernel vector</li> <li><math>C_n = k(x^{(i)}, x^{(j)}) + \sigma^2 I_m</math></li> <li><math>c = k(x^{(n+1)}, x^{(n+1)}) + \sigma^2 I_m</math></li> </ul>
---

Proof:
<ul style="list-style-type: none"> <li>We derive the joint distribution <math>p\left(\begin{bmatrix} y \\ y_{n+1} \end{bmatrix}\right) \sim \mathcal{N}\left[0, \begin{bmatrix} C_n &amp; k \\ k^{\top} &amp; c \end{bmatrix}\right]</math></li> <li>To obtain a closed-form solution for this, we can make use of the following theorem: <ul style="list-style-type: none"> <li>Given a joint Gaussian distribution: <math>\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \sim \mathcal{N}\left[\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} &amp; \Sigma_{12} \\ \Sigma_{21} &amp; \Sigma_{22} \end{bmatrix}\right]</math></li> </ul> </li> </ul>

<ul style="list-style-type: none"> <li>The conditional Gaussian distribution is given by: <math>p(a_2   a_1 = z) = \mathcal{N}(u_2 + \Sigma_{21} \Sigma_{11}^{-1} (z - u_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})</math></li> <li>Then, we get <math>p(y_{n+1}) = \mathcal{N}(k^{\top} C_n^{-1} y, c - k^{\top} C_n^{-1} k)</math></li> </ul>
--

<i>Algorithm</i> —
<ol style="list-style-type: none"> <li>Compute kernel matrix based on observed data</li> <li>Compute kernel vector based on observed data and new instance</li> <li>Calculate mean and variance of predicted distribution</li> <li>Return predicted distribution</li> </ol>

<b>19 PCA</b>
<b>Description</b>
<i>Task</i> — Dimensionality reduction via projection, create uncorrelated features
<i>Description</i> —
<ul style="list-style-type: none"> <li>Unsupervised</li> <li>Non-parametric</li> </ul>

<i>Overview</i> — Identifies lower-dimensional subspace and projects data onto it such that the maximum amount of variance in the data is preserved. In lower-dimensional subspace:
---

<ul style="list-style-type: none"> <li>Axes are called <i>principal components</i>, where the first principal component is the axis accounting for the largest variance</li> <li>Each axis is given by an eigenvector with <i>loadings</i>, indicating how much each variable in the original data contributes to this eigenvector</li> <li>Variance captured along each axis is given by the corresponding eigenvalue</li> </ul>
---

<b>Formulation</b>
<ul style="list-style-type: none"> <li>Project data <math>\{x^{(i)}\}_{i=1}^n \in \mathbb{R}^m</math> onto space <math>\mathbb{R}^d</math> spanned by orthonormal basis <math>\{u^{[j]}\}_{j=1}^d \in \mathbb{R}^m</math> where <math>d \ll m</math></li> <li>Each instance <math>x^{(i)}</math> is projected onto each basis vectors <math>u^{[j]} \cdot x^{(i)}</math></li> <li>Each basis vector <math>u^{[j]}</math> contains <math>m</math> loadings <math>[u_i^{[j]}, \dots, u_m^{[j]}]</math>, whose value indicates how important each feature <math>m</math> is for the <math>j^{th}</math> principal component</li> <li>Mean of projected data for a given basis vector: <math>u^{[j]} \cdot \bar{X} = u^{[j]} \cdot \frac{1}{n} \sum_{i=1}^n x^{(i)}</math></li> <li>Variance of projected data for a given basis vector: <math>\frac{1}{n} \sum_{i=1}^n (u^{[j]} \cdot x^{(i)} - u^{[j]} \cdot \bar{X})^2 = u^{[j]\top} S u^{[j]}</math> where <math>S</math> is the covariance matrix <math>S = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \bar{X})(x^{(i)} - \bar{X})^{\top} = \frac{1}{n} X^{\top} X</math></li> </ul>

<b>Optimization</b>
<i>Parameters</i> — Find $\{u^{[j]}\}_{j=1}^d$
<i>Objective function</i> —
<ul style="list-style-type: none"> <li>Maximize variance <math>\sum_{j=1}^d u^{[j]\top} S u^{[j]}</math> subject to orthonormal <math>\{u^{[j]}\}_{j=1}^d</math></li> <li>Gives rise to Lagrangian formulation</li> <li>Lagrangian formulation for <math>u^{[1]}</math> capturing the most variance: <math>\mathcal{L} = u^{[1]\top} S u^{[1]} - \lambda^{[1]}(u^{[1]} \cdot u^{[1]} - 1)</math> where</li> </ul>



$\lambda^{[1]}$  captures the orthonormality constraint that  $u^{[1]} \cdot u^{[1]} = 1$

- Lagrangian formulation for  $u^{[2]}$  capturing the secondmost variance:  $\mathcal{L} = u^{[2]\top} S u^{[2]} - \lambda^{[2]}(u^{[2]} \cdot u^{[2]} - 1) - \lambda^{[1][2]}(u^{[1]} \cdot u^{[2]} - 0)$  where  $\lambda^{[1][2]}$  captures the orthogonality constraint that  $u^{[1]} \cdot u^{[2]} = 0$

**Optimization** — For  $u^{[1]}$ :

- $\nabla_{u^{[1]}} \mathcal{L} = 2S u^{[1]} - 2\lambda^{[1]} u^{[1]} = 0$
- $\Rightarrow S u^{[1]} = \lambda^{[1]} u^{[1]}$
- This is the eigenvector/eigenvalue equation, so  $u^{[1]}$  is the eigenvector of  $S$  and  $\lambda^{[1]}$  is the associated eigenvalue
- We see that the variance of the projected data is equal to  $\lambda^{[1]}$ :  $u^{[1]\top} S u^{[1]} = u^{[1]\top} \lambda^{[1]} u^{[1]} = \lambda^{[1]} u^{[1]\top} u^{[1]} = \lambda^{[1]} \times 1$

For  $u^{[2]}$ :

- $\nabla_{u^{[2]}} \mathcal{L} = 2S u^{[2]} - 2\lambda^{[2]} u^{[2]} - \lambda^{[1][2]} u^{[1]} = 0$
- $\Rightarrow S u^{[2]} = \lambda^{[2]} u^{[2]}$
- Proof:
  - Multiplying with  $u^{[1]\top}$ :  $2u^{[1]\top} S u^{[2]} - 2\lambda^{[2]} u^{[1]\top} u^{[2]} - \lambda^{[1][2]} u^{[1]\top} u^{[1]} = 0$
  - $= 2u^{[1]\top} S u^{[2]} - 0 - \lambda^{[1][2]} \times 1 = 0$  because of orthogonality resp. orthonormality
  - $= 2u^{[2]\top} S u^{[1]} - \lambda^{[1][2]} = 0$  because the variance is a scalar and can be transposed and because the covariance matrix is symmetric
  - $= 2u^{[2]\top} \lambda^{[1]} u^{[1]} - \lambda^{[1][2]} = 0$  after plugging in the first found basis vector
  - $= 2\lambda^{[1]} \times 0 - \lambda^{[1][2]} = 0$
  - $= \lambda^{[1][2]} = 0$

... continue as for previous vector

In the end, we have a total projected

variance of  $\sum_{j=1}^d \lambda^{[j]}$

**Characteristics** —

- Convex
- Has global minimum
- Has unique or infinitely many solutions
- Can be solved analytically

## 20 GMM

### Description

**Task** — Clustering

**Description** —

- Unsupervised
- Non-parametric

### Formulation

- Instances  $\{x^{(i)}\}_{i=1}^n$
- Each instance has a latent cluster assignment given by:  $z^{(i)} \in \{1, \dots, k\}$
- Probability that cluster assigned to instance  $i$  is cluster  $j$  is given by:  $\pi^{[j]} = p(z^{(i)} = j)$
- Contingent on cluster assignment, each instance is the outcome of a random variable associated with a given cluster:  $x^{(i)}|z^{(i)} = j \sim \mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]})$  where  $\mu^{[j]}$  is the mean and  $\Sigma^{[j]}$  is the covariance associated with cluster  $j$

- Then, marginal distribution of each instance is given by:

$$p(x^{(i)}) = \sum_{j=1}^k \mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]}) \times \pi^{[j]}$$

- This is the GMM, characterized by parameters  $\{\mu^{[j]}, \Sigma^{[j]}, \pi^{[j]}\}_{j=1}^k$

### Optimization

**Parameters** — Find parameters

$\{\mu^{[j]}, \Sigma^{[j]}, \pi^{[j]}\}_{j=1}^k$

**Objective function** —

- Maximize likelihood  $L = \sum_{i=1}^n \log(\sum_{j=1}^k \mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]}) \times \pi^{[j]})$  subject to  $\sum_{j=1}^k \pi^{[j]} = 1$  and  $\Sigma^{[j]} \succ 0$

- This is a constrained, not concave, not analytically solvable optimization problem
- Temporarily assume we know which cluster each instance is associated with
- Let us define a distribution  $q$  over  $1, \dots, k$ :

$$q(z^{(i)}) = p(z^{(i)} = j|x^{(i)}, \theta^{(t)}) = \frac{\mathcal{N}(x^{(i)}|\mu^{[j]}(t), \Sigma^{[j]}(t)) \times \pi^{[j]}(t)}{\sum_{j=1}^k \mathcal{N}(x^{(i)}|\mu^{[j]}(t), \Sigma^{[j]}(t)) \times \pi^{[j]}(t)}$$

- Then, we can rewrite log likelihood as:  $L = \sum_{i=1}^n \log(\sum_{j=1}^k q(z^{(i)}) \frac{\mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]}) \times \pi^{[j]}}{q(z^{(i)})})$

- According to Jensen's inequality:  $L = \sum_{i=1}^n \log(\sum_{j=1}^k q(z^{(i)}) \frac{\mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]}) \times \pi^{[j]}}{q(z^{(i)})}) \geq \sum_{i=1}^n \sum_{j=1}^k q(z^{(i)}) \log(\frac{\mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]}) \times \pi^{[j]}}{q(z^{(i)})})$

- RHS can be rewritten:  $\mathbb{E}_q[\log(p_\theta(x^{(i)}))] =$

$$\mathbb{E}_q[\log(\frac{p_\theta(x^{(i)}, z^{(i)})}{p_\theta(z^{(i)}|x^{(i)})} \frac{q(z^{(i)})}{q(z^{(i)})})] =$$

$$\mathbb{E}_q[\log(\frac{p_\theta(x^{(i)}, z^{(i)})}{q(z^{(i)})})] +$$

$$\mathbb{E}_q[\log(\frac{q(z^{(i)})}{p_\theta(z^{(i)}|x^{(i)})})] = M + E$$

- $E$  corresponds to the KL divergence

between  $q(z^{(i)})$  and  $p(z^{(i)} = j|x^{(i)})$

- $L \geq M \Leftrightarrow E \geq 0$ , which we can show to be the case:

$$- E = \mathbb{E}_q[\log(\frac{q(z^{(i)})}{p_\theta(z^{(i)}|x^{(i)})})] =$$

$$\mathbb{E}_q[-\log(\frac{p_\theta(z^{(i)}|x^{(i)})}{q(z^{(i)})})]$$

- According to Jensen's inequality:

$$E \geq -\log(\mathbb{E}_q[\frac{p_\theta(z^{(i)}|x^{(i)})}{q(z^{(i)})}]) =$$

$$-\log(\sum_{i=1}^k q(z^{(i)}) \frac{p_\theta(z^{(i)}|x^{(i)})}{q(z^{(i)})}) =$$

$$-\log(\sum_{i=1}^k p_\theta(z^{(i)}|x^{(i)}) = -\log(1) = 0$$

- $L = M \Leftrightarrow E = 0$ , i.e. when  $q(z^{(i)}) = p(z^{(i)} = j|x^{(i)}, \theta^{(t)})$
- Then, we have a lower bound on  $L$ , provided by  $M$ , with equality to  $M$ , if we set  $q$  correspondingly
- $M$  is tractable to optimize, since the logarithm now only contains a product, not a sum, and can be decomposed:  $\log(p_\theta(x^{(i)}, z^{(i)})) = \log(\mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]}) \times$

$\pi^{[j]}) = \log(\mathcal{N}(x^{(i)}|\mu^{[j]}, \Sigma^{[j]})) + \log(\pi^{[j]})$

**Optimization** — *Expectation maximization algorithm*

- Randomly initialize  $\theta^{(t)} = \{\mu^{[j]}(t), \Sigma^{[j]}(t), \pi^{[j]}(t)\}_{j=1}^k$
- E-step*: Minimize  $E$ , by computing  $q(z^{(i)})$  given  $x^{(i)}$  and  $\theta^{(t)}$
- M-step*: Maximize  $M$ , by updating  $\theta^{(t)}$  based on MLE for Gaussians, while keeping  $q(z^{(i)})$  fixed:

$$\mu^{[j]}(t+1) = \frac{\sum_{i=1}^n q(z^{(i)}=j) x^{(i)}}{\sum_{i=1}^n q(z^{(i)}=j)}$$

$$\Sigma^{[j]}(t+1) = \frac{\sum_{i=1}^n q(z^{(i)}=j) (x^{(i)} - \mu^{[j]}(t+1))(x^{(i)} - \mu^{[j]}(t+1))^\top}{\sum_{i=1}^n q(z^{(i)}=j)}$$

$$\pi^{[j]}(t+1) = \frac{1}{n} \sum_{i=1}^n q(z^{(i)}=j)$$

- Repeat 2 and 3 until convergence

**Characteristics** —

- Not convex
- May converge to local minimum
- Not analytically solvable
- Always converges, since  $L \geq M$  and  $M^{(t+1)} \geq M^{(t)}$  due to maximizing over  $M$  at each step

## 21 Bayesian Neural Networks

### Setting

- In Bayesian setting, normalization constant is computationally intractable

### Formulation

Since original setting is computationally intractable, we can turn to *variational inference*:

- Variational inference approximates true posterior  $p(w|D)$  by simpler, parametrized distribution  $q(w|\theta)$

- We assume  $q(w|\theta) \sim \mathcal{N}(\mu, \sigma^2 I)$  with  $\theta = (\mu, \sigma)$

### Optimization

**Parameters** — Find parameters  $\theta$

**Objective function** —

- Minimize KL divergence:  $\theta^* = \argmin_{\theta} KL[q(w|\theta)|p(w|D)] = \argmin_{\theta} \mathbb{E}_{w \sim q} \log(q(w|\theta)) - \mathbb{E}_{w \sim q} \log(p(D|w)) - \mathbb{E}_{w \sim q} \log(p(w))$
- Proof:
  - $\argmin_{\theta} KL[q(w|\theta)|p(w|D)] = \argmin_{\theta} \mathbb{E}_{w \sim q} [\log(\frac{q(w|\theta)}{p(w|D)})] = \argmin_{\theta} \mathbb{E}_{w \sim q} [\log(q(w|\theta))] - \mathbb{E}_{w \sim q} \log(p(w|D)] = \argmin_{\theta} \mathbb{E}_{w \sim q} \log(q(w|\theta)) - \mathbb{E}_{w \sim q} \log(p(D|w)) - \mathbb{E}_{w \sim q} \log(p(w)) + \argmin_{\theta} \mathbb{E}_{w \sim q} \log(q(w|\theta)) - \mathbb{E}_{w \sim q} \log(p(D|w)) - \mathbb{E}_{w \sim q} \log(p(w)) + \text{const.}$

**Optimization** —

- To calculate gradient, we can leverage the *reparametrization trick*

- $\frac{\partial}{\partial \theta} \mathbb{E}_{w \sim q} [\log(q(w|\theta)) - \log(p(D|w)) - \log(p(w))] = \frac{\partial}{\partial \theta} \mathbb{E}_{w \sim q} [F(w, \theta)]$  can be reparametrized to:
  - $\frac{\partial}{\partial \mu} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\epsilon^\top \frac{\partial}{\partial w} F(w, \theta) + \frac{\partial}{\partial \mu} F(w, \theta)]$
  - $\frac{\partial}{\partial \sigma} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\epsilon^\top \frac{\partial}{\partial w} F(w, \theta) + \frac{\partial}{\partial \sigma} F(w, \theta)]$
- To optimize this, we can use gradient descent with the following algorithm:
  - Initialize  $\mu$  and  $\sigma$
  - For  $t = 1, 2, \dots$ 
    - Sample  $\epsilon \sim \mathcal{N}(0, I)$
    - Compute  $F(w, \theta)$
    - $\mu_{t+1} \leftarrow \mu_t - \eta_t [\frac{\partial}{\partial w} F(w, \theta) + \frac{\partial}{\partial \mu} F(w, \theta)]_{\mu=\mu_t}$
    - $\sigma_{t+1} \leftarrow \sigma_t - \eta_t [\epsilon^\top \frac{\partial}{\partial w} F(w, \theta) + \frac{\partial}{\partial \sigma} F(w, \theta)]_{\sigma=\sigma_t}$

## 22 Other

### ML Models

**Score** — The score is the derivative of the

log-likelihood:  $\Lambda = \frac{\partial}{\partial \theta} \log(p(x|\theta)) = \frac{\frac{\partial}{\partial \theta} p(x|\theta)}{p(x|\theta)}$

The expected score is given by:

$$\mathbb{E}(\Lambda) = \int p(x|\theta) \frac{\frac{\partial}{\partial \theta} p(x|\theta)}{p(x|\theta)} dx = \frac{\partial}{\partial \theta} \int p(x|\theta) dx = \frac{\partial}{\partial \theta} \times 1 = 0$$

**Fisher information** —

- $I = \mathbb{E}[(\Lambda)^2] = \mathbb{E}[(\frac{\partial}{\partial \theta} \log(p(x|\theta)))^2] = \mathbb{V}(\frac{\partial \log(p(x|\theta))}{\partial \theta})$  where  $\Lambda$  is the score
- Equality is given because  $\mathbb{V}(x) = \mathbb{E}[(x - \mathbb{E}(x))^2] = \mathbb{E}[x^2]$  if  $\mathbb{E}(x) = 0$ , which is the case here, given that the expected score is 0
- Rao-Cramer bound —
  - Shows that there does not exist an asymptotically unbiased parameter estimator
  - For each unbiased estimator,  $\mathbb{E}[(\hat{\theta} - \theta)^2] \geq \frac{1}{I}$  where  $I$  is the Fisher information
  - For estimators in general,

$$(\frac{\frac{\partial}{\partial \theta} \text{bias} + 1}{I} + \text{bias}^2 \leq \mathbb{E}[(\hat{\theta} - \theta)^2], \text{ so there is a trade-off if the bias derivative is negative and the squared bias is positive, whereby a biased estimator may produce better results than an unbiased estimator}$$

is a trade-off if the bias derivative is negative and the squared bias is positive, whereby a biased estimator may produce better results than an unbiased estimator

**Proof:**

- Given Cauchy Schwarz inequality, we can say:  $\mathbb{E}[(\Lambda - \mathbb{E}(\Lambda))(\hat{\theta} - \mathbb{E}(\hat{\theta}))]^2 \leq \mathbb{E}[(\Lambda - \mathbb{E}(\Lambda))^2] \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2]$  where  $\Lambda$  is the score
- We know that  $\mathbb{E}(\Lambda) = 0$
- Let's look at the LHS of the equation:
  - Since  $\mathbb{E}(\Lambda) = 0$ , we can simplify to  $\mathbb{E}[\Lambda(\hat{\theta} - \mathbb{E}(\hat{\theta}))] = \mathbb{E}[\Lambda \hat{\theta}] - \mathbb{E}[\Lambda] \mathbb{E}[\hat{\theta}] = \mathbb{E}[\Lambda \hat{\theta}] - 0$
  - This can be developed to:

$$\mathbb{E}[\Lambda \hat{\theta}] = \int p(x|\theta) \frac{\frac{\partial}{\partial \theta} p(x|\theta)}{p(x|\theta)} \hat{\theta} dx =$$

$$\frac{\partial}{\partial \theta} (\int p(x|\theta) \hat{\theta} dx - \theta) + 1 \text{ where the last part } (-\theta + 1) \text{ can be added, because}$$

- $\frac{\partial}{\partial \theta} - \theta = -1$  and we compensate this with +1
  - This is equal to the derivative of the bias + 1:  $\frac{\partial}{\partial \theta} (\int p(x|\theta) \hat{\theta} dx - \theta) + 1 = \frac{\partial}{\partial \theta} (\mathbb{E}[\hat{\theta}] - \theta) + 1 = \frac{\partial}{\partial \theta} \text{bias} + 1$
- Let's look at the RHS of the equation: Since  $\mathbb{E}(\Lambda) = 0$ , first term is  $\mathbb{E}(\Lambda^2) = I$
- Then, we have:  $(\frac{\partial}{\partial \theta} \text{bias} + 1)^2 \leq I \times \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2] = I \times \mathbb{E}[(\hat{\theta} - \theta - \mathbb{E}(\hat{\theta}) + \theta)^2] = \dots = I \times \mathbb{E}[(\hat{\theta} - \theta)^2] - \text{bias}^2$
- Then, we have  $(\frac{\frac{\partial}{\partial \theta} \text{bias} + 1}{I} + \text{bias}^2 \leq \mathbb{E}[(\hat{\theta} - \theta)^2]$

### Causal Models

**Causal scenarios** —

- Causal scenario without selection bias:  $\mathcal{X}$  affects  $\mathcal{Y}$  and there is no selection bias
  - Some features  $\mathcal{X}_{\perp Y}$  do not causally affect  $\mathcal{Y}$ , but are affected by  $\mathcal{W}$
  - Some features  $\mathcal{X}_{\perp W}$  causally affect  $\mathcal{Y}$ , but are not affected by  $\mathcal{W}$
  - Some features  $\mathcal{X}_{W \& Y}$  causally affect  $\mathcal{Y}$  and are affected by  $\mathcal{W}$  as well as  $\mathcal{X}_{\perp Y}$  and  $\mathcal{X}_{\perp W}$

- Anti causal scenario: We assume  $\mathcal{Y}$  affects  $\mathcal{X}$ , rather than the other way around

- Causal scenario with selection bias:  $\mathcal{X}$  affects  $\mathcal{Y}$  and there is a selection bias

**Counterfactual invariance** —

- Counterfactual invariance: Results of estimator remain consistent across different counterfactual scenarios, i.e. if  $\mathcal{Y}$  is affected by  $\mathcal{X}$ , and  $\mathcal{X}$  is affected by  $\mathcal{W}$ , but  $\mathcal{W}$  does not affect  $\mathcal{Y}$ , our estimator should be invariant to states of  $\mathcal{W}$ , i.e.  $f(\mathcal{X}(\mathcal{W}_1)) = f(\mathcal{X}(\mathcal{W}_2))$
- For counterfactual invariance, the following must hold:
  - Causal scenario without selection bias:  $f(\mathcal{X}_{\perp W})|Y$ , i.e. estimate  $f$  only depends on  $\mathcal{X}_{\perp W}$
  - Anti causal scenario:  $(f(\mathcal{X}) \perp W)|Y$ , i.e. estimate  $f$  only depends on  $\mathcal{X}_{\perp W}$ , provided  $\mathcal{Y}$  is known
  - Causal scenario with selection bias:  $(f(\mathcal{X}) \perp W)|Y$  as long as  $\mathcal{X}_{\perp Y}$  and  $\mathcal{X}_{W \& Y}$  do not influence  $\mathcal{Y}$  whatsoever, i.e.  $(Y \perp \mathcal{X})|X_{\perp W}, W$

- For causal scenario without selection bias we need to show:  $\mathcal{X}_{\perp W} \perp W$

- For anti causal scenario we need to show:  $(\mathcal{X}_{\perp W} \perp W)|Y$

- This can be shown via *d-separation*

**D separation** —

- Undirected path of  $n$  nodes is d-separated, if it contains 3 nodes following any of the following forms and if this form is blocked:
  - Chain structure:  $X \rightarrow Z \rightarrow Y$  or  $Y \rightarrow Z \rightarrow X$  — is blocked, if we condition on  $Z$ , i.e.  $Z$  is known
  - Fork structure:  $X \leftarrow Z \rightarrow Y$  — is blocked, if we condition on  $Z$ , i.e.  $Z$  is known
  - Collider structure:  $X \rightarrow Z \leftarrow Y$  — is blocked, if we don't condition on  $Z$  or any of its descendants
- Random variables  $X$  and  $Y$  are conditionally independent if each path

between them is d-separated

→ as soon as we have one blocked triple

on path, entire path is blocked

→ as soon as one path is active, we cannot

guarantee conditional independence

- For causal scenario without selection bias we can show  $\mathcal{X} \perp_W \perp \mathcal{W}$  since all paths are blocked
- For anti causal scenario we can show  $(\mathcal{X} \perp_W \perp \mathcal{W}) | \mathcal{Y}$  since all paths are blocked, conditioned on  $\mathcal{Y}$ , i.e. if  $\mathcal{Y}$  is observed