# *Statistical Machine Learning*

Christian Walder

Machine Learning Research Group
CSIRO Data61

and

College of Engineering and Computer Science
The Australian National University

Canberra
Semester One, 2020.

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")

Part II

*Introduction*

# *Flavour of this course*

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

- Formalise intuitions about problems
- Use language of mathematics to express models
- Geometry, vectors, linear algebra for reasoning
- Probabilistic models to capture uncertainty
- Design and analysis of algorithms
- Numerical algorithms in python
- Understand the choices when designing machine learning methods

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*
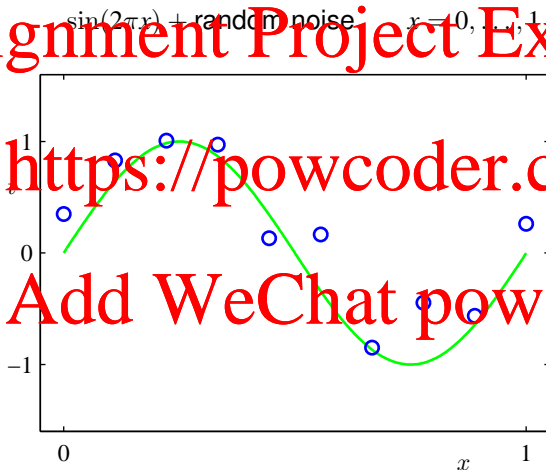
*Expectations and Covariances*

**Definition (Mitchell, 1998)**

A computer program is said to learn from experience $E$ with respect to some class of tasks $T$ and performance measure $P$, if its performance at tasks in $T$, as measured by $P$, improves with experience $E$.

# Polynomial Curve Fitting

- some artificial data created from the function

$$\sin(2\pi x) + \text{random noise} \quad x = 0, \ldots, 1$$

$$N = 10$$
$$\mathbf{x} \equiv (x_1, \ldots, x_N)^T$$
$$\mathbf{t} \equiv (t_1, \ldots, t_N)^T$$

*Polynomial Curve Fitting*

*Probability Theory*
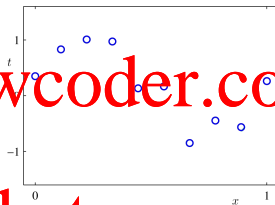
*Probability Densities*

*Expectations and Covariances*

$N = 10$

$\mathbf{x} \equiv (x_1, \ldots, x_N)^T$

$\mathbf{t} \equiv (t_1, \ldots, t_N)^T$

$x_i \in \mathbb{R} \quad i = 1, \ldots, N$

$t_i \in \mathbb{R} \quad i = 1, \ldots, N$

# Polynomial Curve Fitting - Model Specification

$M$: order of polynomial

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \cdots + w_M x^M$$

$$= \sum_{m=0}^{M} w_m x^m$$

- nonlinear function of $x$
- linear function of the unknown model parameter $\mathbf{w}$
- How can we find good parameters $\mathbf{w} = (w_1, \ldots, w_M)^T$?

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

# *Learning is Improving Performance*

- Performance measure : Error between target and prediction of the model for the training data

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (y(x_n, \mathbf{w}) - t_n)^2$$

- unique minimum of $E(\mathbf{w})$ for argument $\mathbf{w}^\star$ under certain conditions (what are they?)

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

$$y(x, \mathbf{w}) = \sum_{m=0}^{M} w_m x^m \Big|_{M=0}$$

$$= w_0$$



$M = 0$

$$y(x, w) \equiv \left. \sum_{m=0}^{M} w_m x^m \right|_{M=1}$$

$$= w_0 + w_1 x$$

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

$$y(x, \mathbf{w}) = \sum_{m=0}^{M} w_m x^m \bigg|_{M=3}$$

$$= w_0 + w_1 x + w_2 x^2 + w_3 x^3$$



$M = 3$

$$y(x, \mathbf{w}) = \sum_{m=0}^{M} w_m\, x^m$$

$$\Big|_{M=9} = w_0 + w_1\, x + \cdots + w_8\, x^8 + w_9\, x^9$$

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

- overfitting



$M = 9$

- Train the model and get $\mathbf{w}^\star$
- Get 100 new data points
- Root mean-square (RMS) error

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^\star)/N}$$

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

| | M = 0 | M = 1 | M = 3 | M = 9 |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ | | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ | | | -25.43 | -5321.83 |
| $w_3^\star$ | | | 17.37 | 48568.31 |
| $w_4^\star$ | | | | -231639.30 |
| $w_5^\star$ | | | | 640042.26 |
| $w_6^\star$ | | | | -1061800.52 |
| $w_7^\star$ | | | | 1042400.18 |
| $w_8^\star$ | | | | -557682.99 |
| $w_9^\star$ | | | | 125201.43 |

*Table:* Coefficients $\mathbf{w}^\star$ for polynomials of various order.

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

$N = 15$

Statistical Machine
Learning

©2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

# *More Data*

- $N = 100$
- heuristics : have no less than 5 to 10 times as many data points than parameters
- but number of parameters is not necessarily the most appropriate measure of model complexity !
- later: Bayesian approach

Assignment Project Exam Help

- How to constrain the growing of the coefficients $\mathbf{w}$?
- Add a regularisation term to the error function

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2}\sum_{n=1}^{N}(y(x_n,\mathbf{w}) - t_n)^2 + \frac{\lambda}{2}\|\mathbf{w}\|^2$$

https://powcoder.com

- Squared norm of the parameter vector $\mathbf{w}$

Add WeChat powcoder

$$\|\mathbf{w}\|^2 = \mathbf{w}^T\mathbf{w} = w_0^2 + w_1^2 + \cdots + w_M^2$$

- unique minimum of $E(\mathbf{w})$ for argument $\mathbf{w}^\star$ under certain conditions (what are they for $\lambda = 0$? for $\lambda > 0$?)

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

• $M = 9$



$\ln \lambda = -18$

- $M = 9$



$\ln \lambda = 0$

Assignment Project Exam Help

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

https://powcoder.com



Add WeChat powcoder

- $M = 9$

# *What is Machine Learning?*

**Definition (Mitchell, 1998)**
A computer program is said to learn from experience $E$ with respect to some class of tasks $T$ and performance measure $P$, if its performance at tasks in $T$, as measured by $P$, improves with experience $E$.

- Task: regression
- Experience: $\mathbf{x}$ input examples, $\mathbf{t}$ output labels
- Performance: squared error
- Model choice
- Regularisation
- **do not train on the test set!**

Polynomial Curve Fitting

**Probability Theory**

Probability Densities

Expectations and Covariances



$p(X,Y)$

$Y=2$

$Y=1$

$X$

*Polynomial Curve Fitting*

**Probability Theory**

*Probability Densities*

*Expectations and Covariances*

| Y vs. X | a | b | c | d | e | f | g | h | i | sum |
|---------|---|---|---|---|---|---|---|---|---|-----|
| 2 | 0 | 0 | 0 | 1 | 4 | 5 | 8 | 6 | 2 | 26 |
| 1 | 3 | 6 | 8 | 8 | 5 | 3 | 0 | 0 | 0 | 34 |
| sum | 3 | 6 | 8 | 9 | 9 | 8 | 8 | 6 | 2 | 60 |

*Polynomial Curve Fitting*

**Probability Theory**

*Probability Densities*

*Expectations and Covariances*

| Y vs. X | a | b | c | d | e | f | g | h | i | sum |
|---------|---|---|---|---|---|---|---|---|---|-----|
| 2 | 0 | 0 | 0 | 1 | 4 | 5 | 8 | 6 | 2 | 26 |
| 1 | 3 | 6 | 8 | 8 | 5 | 3 | 1 | 0 | 0 | 34 |
| sum | 3 | 6 | 8 | 9 | 9 | 8 | 9 | 6 | 2 | 60 |

$$p(X = d, Y = 1) = 8/60$$
$$p(X = d) = p(X = d, Y = 2) + p(X = d, Y = 1)$$
$$= 1/60 + 8/60$$

$$p(X = d) = \sum_Y p(X = d, Y)$$

$$p(X) = \sum_Y p(X, Y)$$

# Sum Rule

Statistical Machine
Learning

©2020
Ong & Walder & Webers
Data61 \ CSIRO
The Australian National
University

Polynomial Curve Fitting

**Probability Theory**

Probability Densities

Expectations and
Covariances

| Y vs. X | a | b | c | d | e | f | g | h | i | sum |
|---------|---|---|---|---|---|---|---|---|---|-----|
| | | | | 1 | 4 | 5 | | 6 | 2 | 26 |
| | 3 | 6 | 8 | 8 | 5 | 3 | 1 | 0 | 0 | 34 |
| sum | 3 | 6 | 8 | 9 | 9 | 8 | 9 | 6 | 2 | 60 |

$$p(X) = \sum_Y p(X, Y) \qquad p(Y) = \sum_X p(X, Y)$$

| Y vs. X | a | b | c | d | e | f | g | h | i | sum |
|---------|---|---|---|---|---|---|---|---|---|-----|
| 2 | 0 | 0 | 0 | 1 | 4 | 5 | 8 | 6 | 2 | 26 |
| 1 | 3 | 6 | 8 | 8 | 5 | 3 | 1 | 0 | 0 | 34 |
| sum | 3 | 6 | 8 | 9 | 8 | 9 | 6 | 2 | | 60 |

**Conditional Probability**

$$p(X = d \mid Y = 1) = 8/34$$

Calculate $p(Y = 1)$:

$$p(Y = 1) = \sum_X p(X, Y = 1) = 34/60$$

$$p(X = d, Y = 1) = p(X = d \mid Y = 1)\, p(Y = 1)$$

$$p(X, Y) = p(X \mid Y)\, p(Y)$$

Another intuitive view is renormalisation of relative frequencies:

$$p(X \mid Y) = \frac{p(X, Y)}{p(Y)}$$

# *Sum and Product Rules*

| Y vs. X | a | b | c | d | e | f | g | h | i | sum |
|---------|---|---|---|---|---|---|---|---|---|-----|
| | 0 | 0 | 0 | 1 | 4 | 5 | 8 | 6 | 2 | 26 |
| | 3 | 6 | 8 | 8 | 5 | 3 | 1 | 0 | 0 | 34 |
| sum | 3 | 6 | 8 | 9 | 9 | 8 | 9 | 6 | 2 | 60 |

$$p(X) = \sum_Y p(X, Y) \qquad p(X \mid Y) = \frac{p(X, Y)}{p(Y)}$$

*Polynomial Curve Fitting*

**Probability Theory**

*Probability Densities*

*Expectations and Covariances*

- Sum Rule

$$p(X) = \sum_Y p(X, Y)$$

- Product Rule

$$p(X, Y) = p(X \mid Y)\, p(Y)$$

These rules form the basis of Bayesian machine learning, and this course!

Polynomial Curve Fitting

**Probability Theory**

Probability Densities

Expectations and Covariances

Use product rule

$$p(X, Y) = p(X \mid Y) p(Y) = p(Y \mid X) p(X)$$

Bayes Theorem

$$p(Y \mid X) = \frac{p(X \mid Y) p(Y)}{p(X)} \qquad \text{only defined for } p(X) > 0$$

and

$$p(X) = \sum_Y p(X, Y) \qquad \text{(sum rule)}$$

$$= \sum_Y p(X \mid Y) p(Y) \qquad \text{(product rule)}$$

- Real valued variable $x \in \mathbb{R}$
- Probability of $x$ to fall in the interval $(x, x + \delta x)$ is given by $p(x)\delta x$ for infinitesimal small $\delta x$.

-
$$p(x \in (a, b)) = \int_a^b p(x) \, dx.$$

- Nonnegative

Assignment Project Exam Help

$$p(x) \geq 0$$

- Normalisation

$$\int_{-\infty}^{\infty} p(x) \, dx = 1.$$

https://powcoder.com

Add WeChat powcoder

DATA 61 · CSIRO

*Polynomial Curve Fitting*

*Probability Theory*

**Probability Densities**

*Expectations and Covariances*

$$P(x) = \int_{-\infty}^{x} p(x) \, dx$$

or

$$\frac{d}{dx} P(x) = p(x)$$

- Vector $\mathbf{x} = (x_1, \ldots, x_D)^T = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}$

- Nonnegative
$$p(\mathbf{x}) \geq 0$$

- Normalisation
$$\int_{-\infty}^{\infty} p(\mathbf{x}) \, d\mathbf{x} = 1.$$

- This means
$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} p(\mathbf{x}) \, dx_1 \ldots \, dx_D = 1.$$

*Polynomial Curve Fitting*

*Probability Theory*

**Probability Densities**

*Expectations and Covariances*

- Sum Rule

$$p(x) = \int_{-\infty}^{\infty} p(x, y) \, dy$$

- Product Rule

$$p(x, y) = p(y \mid x) \, p(x)$$

*Polynomial Curve Fitting*

*Probability Theory*

*Probability Densities*

*Expectations and Covariances*

- Weighted average of a function f(x) under the probability distribution $p(x)$

$$\mathbb{E}[f] = \sum_x p(x) f(x) \qquad \text{discrete distribution } p(x)$$

$$\mathbb{E}[f] = \int p(x) f(x) \, dx \qquad \text{probability density } p(x)$$

- Given a finite number $N$ of points $x_n$ drawn from the probability distribution $p(x)$.

- Approximate the expectation by a finite sum:

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^{N} f(x_n)$$

- How to draw points from a probability distribution $p_n(x)$?
  Lecture coming about "Sampling"

Polynomial Curve Fitting

Probability Theory

Probability Densities

*Expectations and*
*Covariances*

- arbitrary function $f(x, y)$

$$\mathbb{E}_x [f(x, y)] = \sum_x p(x) f(x, y) \qquad \text{discrete distribution } p(x)$$

$$\mathbb{E}_x [f(x, y)] = \int p(x) f(x, y) \, dx \qquad \text{probability density } p(x)$$

- Note that $\mathbb{E}_x [f(x, y)]$ is a function of $y$.

- arbitrary function $f(x)$

$$\mathbb{E}_x[f \mid y] = \sum p(x \mid y)f(x) \qquad \text{discrete distribution } p(x)$$

$$\mathbb{E}_x[f \mid y] = \int p(x \mid y)f(x)\, dx \qquad \text{probability density } p(x)$$

- Note that $\mathbb{E}_x[f \mid y]$ is a function of $y$.
- Other notation used in the literature $\mathbb{E}_{x|y}[f]$.
- What is $\mathbb{E}[\mathbb{E}[f(x) \mid y]]$? Can we simplify it?
- This must mean $\mathbb{E}_y[\mathbb{E}_x[f(x) \mid y]]$. (Why?)

$$\mathbb{E}_y[\mathbb{E}_x[f(x) \mid y]] = \sum_y p(y)\,\mathbb{E}_x[f \mid y] = \sum_y p(y)\sum_x p(x \mid y)f(x)$$

$$= \sum_{x,y} f(x)\,p(x,y) = \sum_x f(x)\,p(x)$$

$$= \mathbb{E}_x[f(x)]$$

- arbitrary function $f(x)$

$$\text{var}[f] = \mathbb{E}\left[(f(x) - \mathbb{E}[f(x)])^2\right] = \mathbb{E}\left[f(x)^2\right] - \mathbb{E}[f(x)]^2$$

- Special case: $f(x) = x$

$$\text{var}[x] = \mathbb{E}\left[(x - \mathbb{E}[x])^2\right] = \mathbb{E}\left[x^2\right] - \mathbb{E}[x]^2$$

# *Covariance*

- Two random variables $x \in \mathbb{R}$ and $y \in \mathbb{R}$

$$\mathrm{cov}[x, y] = \mathbb{E}_{x,y}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])]$$
$$= \mathbb{E}_{x,y}[x\,y] - \mathbb{E}[x]\,\mathbb{E}[y]$$

- With $\mathbb{E}[x] = a$ and $\mathbb{E}[y] = b$

$$\mathrm{cov}[x, y] = \mathbb{E}_{x,y}[(x - a)(y - b)]$$
$$= \mathbb{E}_{x,y}[x\,y] - \mathbb{E}_{x,y}[x\,b] - \mathbb{E}_{x,y}[a\,y] + \mathbb{E}_{x,y}[a\,b]$$
$$= \mathbb{E}_{x,y}[x\,y] - b\,\underbrace{\mathbb{E}_{x,y}[x]}_{=\mathbb{E}[x]} - a\,\underbrace{\mathbb{E}_{x,y}[y]}_{=\mathbb{E}[y]} + a\,b\,\underbrace{\mathbb{E}_{x,y}[1]}_{=1}$$
$$= \mathbb{E}_{x,y}[x\,y] - a\,b - a\,b + a\,b = \mathbb{E}_{x,y}[x\,y] - a\,b$$
$$= \mathbb{E}_{x,y}[x\,y] - \mathbb{E}[x]\,\mathbb{E}[y]$$

- Expresses how strongly $x$ and $y$ vary together. If $x$ and $y$ are independent, their covariance vanishes.

Assignment Project Exam Help

https://powcoder.com

- Two random variables $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^D$

$$\text{cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y}^{\mathbf{T}} - \mathbb{E}[\mathbf{y}^{\mathbf{T}}]) \right]$$
$$= \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ \mathbf{x}\,\mathbf{y}^{\mathbf{T}} \right] - \mathbb{E}[\mathbf{x}]\,\mathbb{E}\left[\mathbf{y}^{\mathbf{T}}\right]$$
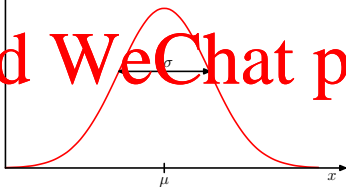
Add WeChat powcoder

- $x \in \mathbb{R}$
- Gaussian Distribution with mean $\mu$ and variance $\sigma^2$

$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\{-\frac{1}{2\sigma^2}(x - \mu)^2\}$$

- $\mathcal{N}(x \mid \mu, \sigma^2) > 0$
- $\int_{-\infty}^{\infty} \mathcal{N}(x \mid \mu, \sigma^2) \, dx = 1$
- Expectation over $x$

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x \mid \mu, \sigma^2) \, x \, dx = \mu$$

- Expectation over $x^2$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x \mid \mu, \sigma^2) \, x^2 \, dx = \mu^2 + \sigma^2$$

- Variance of x

$$\mathrm{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

Assignment Project Exam Help

- Estimate best predictor = training = learning
  Given data $(x_1, y_1), \ldots, (x_n, y_n)$, find a predictor $f_{\mathbf{w}}(\cdot)$.

  https://powcoder.com

  1. Identify the type of input $x$ and output $y$ data
  2. Choose a linear mathematical model for $f_{\mathbf{w}}$
  3. Design an objective function or likelihood
  4. Calculate the optimal parameter ($\mathbf{w}$)
  5. Model uncertainty using the Bayesian approach
  6. Implement and compute (the algorithm in python)
  7. Interpret and diagnose results

Add WeChat powcoder