

## 1. SIMPLE LINEAR REGRESSION

1.1. **The regression line.** Consider a data set

$$\mathcal{D} = \{(x_i, y_i) : i = 1, \dots, n\}.$$

If the *mean-squared error* function

$$\text{MSE}(a, b) = \frac{1}{n} \sum_{i=1}^n (ax_i + b - y_i)^2$$

achieves its absolute minimum value at

$$(a, b) = (\alpha, \beta)$$

then the line  $y = \alpha x + \beta$  is called the *regression line* or *least-squares line* for  $\mathcal{D}$ .

The *slope*,  $\alpha$ , and the *intercept*,  $\beta$  of the regression line (its *coefficients*) can be expressed in terms of basic statistics of  $\mathcal{D}$ :

means:	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i,$	$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$
variances:	$s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2,$	$s_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$
covariance:	$s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$	

**Theorem 1** (Gauss/Legendre). *The coefficients of the regression line of  $\mathcal{D}$  are:*

$$a = \frac{s_{xy}}{s_x^2}, \quad b = \bar{y} - a\bar{x}.$$

*Proof.* Notice that

$$\min_{(a,b)} \text{MSE}(a, b) = \min_a \left( \min_b \text{MSE}(a, b) \right).$$

For a given  $a$ , the quantity  $\text{MSE}(a, b)$  is a quadratic polynomial in  $b$ :

$$\text{MSE}(a, b) = b^2 - 2 \left( \frac{1}{n} \sum_{i=1}^n (y_i - ax_i) \right) b + \sum_{i=1}^n (y_i - ax_i)^2$$

Since a quadratic polynomial  $t^2 - 2qt + r$  achieves its minimum value at  $t = q$ ,  $\text{MSE}(a, b)$  achieves its minimum value when

$$b = \frac{1}{n} \sum_{i=1}^n (y_i - ax_i) = \bar{y} - a\bar{x}.$$

It remains to determine

$$\min_a \text{MSE}(a, \bar{y} - a\bar{x}) = \min_a \frac{1}{n} \sum_{i=1}^n (ax_i + (\bar{y} - a\bar{x}) - y_i)^2.$$

Expanding and rearranging, we get

$$\frac{1}{n} \sum_{i=1}^n (ax_i + (\bar{y} - a\bar{x}) - y_i)^2 = s_x^2 a^2 - 2s_{xy}a + s_y^2.$$

Since a quadratic polynomial  $pt^2 - 2qt + r$  achieves its minimum value at  $t = q/p$ , the function  $\text{MSE}(a, \bar{y} - a\bar{x})$  achieves its minimum value when  $a = s_{xy}/s_x^2$ .

Thus,  $\text{MSE}(a, b)$  is minimized when

$$a = \frac{s_{xy}}{s_x^2}, \quad b = \bar{y} - a\bar{x}. \quad \square$$

Define  $\mathbf{1}, \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  by

$$\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

For  $\alpha, \beta \in \mathbb{R}$ , define the associated *residual vector*,  $\mathbf{e}(\alpha, \beta)$ , by

$$\mathbf{e}(\alpha, \beta) = \alpha\mathbf{x} + \beta\mathbf{1} - \mathbf{y}.$$

Then

$$\text{MSE}(\alpha, \beta) = \frac{1}{n} \|\mathbf{e}(\alpha, \beta)\|^2.$$

Let  $U$  be the subspace of  $\mathbb{R}^n$  spanned by the vectors  $\mathbf{x}$  and  $\mathbf{1}$ :

$$U = \{\alpha\mathbf{x} + \beta\mathbf{1} : \alpha, \beta \in \mathbb{R}\}.$$

Let  $d(\mathbf{y}, U)$  be the distance from  $\mathbf{y}$  to  $U$ , i.e., the minimal distance from  $\mathbf{y}$  to an element of  $U$ :

$$d(\mathbf{y}, U) = \inf_{a,b} \|a\mathbf{x} + b\mathbf{1} - \mathbf{y}\|.$$

The infimum on the right is achieved by *orthogonal projection of  $\mathbf{y}$  onto  $U$* , i.e., the unique vector  $\hat{\mathbf{y}} \in U$  such that

$$\langle \hat{\mathbf{y}}, \mathbf{y} - \hat{\mathbf{y}} \rangle = 0.$$

If  $\{\mathbf{u}_1, \mathbf{u}_2\}$  is any orthonormal basis of  $U$ , then

$$\hat{\mathbf{y}} = \langle \mathbf{u}_1, \mathbf{y} \rangle \mathbf{u}_1 + \langle \mathbf{u}_2, \mathbf{y} \rangle \mathbf{u}_2.$$

We can construct an orthonormal basis of  $U$  by applying the *Gram-Schmidt orthonormalization procedure* to the spanning set  $\{\mathbf{1}, \mathbf{x}\}$ . Let

$$\begin{aligned} \mathbf{u}_1 &= \frac{1}{\|\mathbf{1}\|} \mathbf{1} = \frac{1}{\sqrt{n}} \mathbf{1}, \\ \mathbf{u}_2' &= \mathbf{x} - \langle \mathbf{u}_1, \mathbf{x} \rangle \mathbf{u}_1 \\ &= \mathbf{x} - \frac{1}{\sqrt{n}} \langle \mathbf{1}, \mathbf{x} \rangle \frac{1}{\sqrt{n}} \mathbf{1} \\ &= \mathbf{x} - \bar{x} \mathbf{1}, \end{aligned}$$

Assume that  $\mathbf{x}$  and  $\mathbf{1}$  are linearly independent. Then  $\mathbf{u}'_2 \neq 0$  and we may set

$$\begin{aligned}\mathbf{u}_2 &= \frac{1}{\|\mathbf{u}'_2\|} \mathbf{u}'_2 \\ &= \frac{1}{\sqrt{n}s_x} (\mathbf{x} - \bar{x}\mathbf{1})\end{aligned}$$

Thus, if  $\mathbf{x}$  and  $\mathbf{1}$  are linearly independent, then

$$\left\{ \frac{1}{\sqrt{n}} \mathbf{1}, \frac{1}{\sqrt{n}s_x} (\mathbf{x} - \bar{x}\mathbf{1}) \right\}.$$

is an orthonormal basis of  $U$ . It follows that

$$\hat{\mathbf{y}} = \frac{1}{n} \langle \mathbf{1}, \mathbf{y} \rangle \mathbf{1} + \frac{1}{ns_x^2} \langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} \rangle (\mathbf{x} - \bar{x}\mathbf{1})$$

Since  $\mathbf{x} - \bar{x}\mathbf{1}$  is orthogonal to  $\mathbf{1}$ ,

$$\frac{1}{n} \langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} \rangle = \frac{1}{n} \langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} - \bar{y}\mathbf{1} \rangle = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = s_{xy}.$$

$$\hat{\mathbf{y}} = \bar{y}\mathbf{1} + \frac{s_{xy}}{s_x^2} (\mathbf{x} - \bar{x}\mathbf{1}) = \frac{s_{xy}}{s_x^2} \mathbf{x} + \left( \bar{y} - \frac{s_{xy}}{s_x^2} \bar{x} \right) \mathbf{1}$$

**Theorem 2.**

(1) *There is a unique vector  $\hat{\mathbf{y}} \in U$  such that*

$$d(\mathbf{y}, U) = \|\mathbf{y} - \hat{\mathbf{y}}\|.$$

(2) *If the vectors  $\mathbf{1}$  and  $\mathbf{x}$  are linearly independent, then there are unique scalars  $\hat{a}$  and  $\hat{b}$  such that*

$$\hat{\mathbf{y}} = \hat{a}\mathbf{x} + \hat{b}\mathbf{1}.$$

$$\|\mathbf{y} - \bar{y}\mathbf{1}\|^2 = \|(\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \bar{y}\mathbf{1})\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \bar{y}\mathbf{1}\|^2 = \text{SSE} + s_{\hat{\mathbf{y}}}^2$$

**1.2. Sums of squares.** The regression line gives the estimate

$$\hat{y}_i = ax_i + b$$

for  $y_i$ . The  $\hat{y}_i$  and the  $y_i$  have the same mean:

$$\bar{\hat{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (ax_i + b) = a\bar{x} + b = \bar{y},$$

the final equality following from Theorem 1.

$$\begin{aligned}
s_y^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\
&= \text{MSE}(a, b) + 2s_{e\hat{y}} + s_{\hat{y}}^2.
\end{aligned}$$

## 2. THE BIVARIATE NORMAL DISTRIBUTION

The bivariate normal density with means  $\mu_1$  and  $\mu_2$ , variances  $\sigma_1$  and  $\sigma_2$ , and correlation  $\rho$  is defined by

$$f(x_1, x_2) = \frac{1}{2\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{1}{2}Q(x_1, x_2)},$$

where

$$Q(x_1, x_2) = \frac{1}{\sqrt{1-\rho^2}} \left[ \left( \frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \left( \frac{x_1 - \mu_1}{\sigma_1} \right) \left( \frac{x_2 - \mu_2}{\sigma_2} \right) + \left( \frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right]$$

We write

$$(X_1, X_2) \sim N(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$$

if  $(X_1, X_2)$  has density  $f(x_1, x_2)$ .

Suppose  $X \sim N(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$ . Prove:

- (1) The marginal density of  $X_1$  is the univariate normal density with mean  $\mu_1$  and variance  $\sigma_1^2$ , i.e.,

$$\int_{-\infty}^{\infty} f(x_1, x_2) dx_2 = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2}\left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2}.$$

- (2)  $E[X_i] = \mu_i$ ,  $E[(X_i - \mu_i)^2] = \sigma_i^2$ , and  $E[(X_1 - \mu_1)(X_2 - \mu_2)] = \sigma_1\sigma_2\rho$ .

- (3) The conditional density of  $X_2$  given  $X_1$  is given by

$$f(x_2|x_1) = \frac{1}{\sqrt{2\pi(1-\rho^2)}\sigma_2} e^{-\frac{1}{2}\left(\frac{x_2 - \left(\rho\frac{\sigma_2}{\sigma_1}(x_1 - \mu_1) + \mu_2\right)}{\sqrt{1-\rho^2}\sigma_2}\right)^2}.$$

- (4) The conditional expectation and variance of  $X_2$  given  $X_1$  are given by

$$E[X_2|X_1] = \rho\frac{\sigma_2}{\sigma_1}(X_1 - \mu_1) + \mu_2$$

and

$$E[(X_2 - E[X_2|X_1])^2|X_1] = \sqrt{1-\rho^2}\sigma_2,$$

respectively. Note that the latter quantity is independent of  $X_1$ .

### 3. CONDITIONAL EXPECTATION

**Theorem-Definition 3.** Let  $\Omega$  be a set equipped with a probability measure,  $P$ . Given random variables  $X$  and  $Y$  on  $\Omega$ , there is a unique function  $f : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$\int_{[X \in G]} Y dP = \int_{[X \in G]} f(X) dP,$$

for every  $E \subseteq \mathbb{R}$ . The random variable  $f(X)$  is called the conditional expectation of  $Y$  given  $X$  and denoted  $E[Y|X]$ .

(1) If  $Y = f(X)$ , then  $E[Y|X] = Y$ .

(2) If  $X = 1$ , then  $E[Y|X] = E[Y]$ :

$$1 \notin G : \quad \int_{[X \in G]} Y dP = \int_{\emptyset} Y dP = 0 = \int_{\emptyset} E[Y] dP = \int_{[X \in G]} E[Y] dP$$

$$1 \in G : \quad \int_{[X \in G]} Y dP = \int_{\Omega} Y dP = E[Y] = \int_{\Omega} E[Y] dP = \int_{[X \in G]} E[Y] dP$$

(3) If  $E[Y|X] = f(X)$ , then

$$E[I_H(X)Y|X] = I_H(X)f(X)$$

for all  $H \subseteq \mathbb{R}$ :

$$\int_{[X \in G]} I_H(X)Y dP = \int_{[X \in G \cap H]} Y dP = \int_{[X \in G \cap H]} f(X) dP = \int_{[X \in G]} I_H(X)f(X) dP$$

(4) If  $u : \mathbb{R} \rightarrow \mathbb{R}$ , then

$$E[u(X)Y|X] = u(X) E[Y|X].$$

(Proof: Exercise?)

(5) If  $X = u(Y)$ , then

$$E[E[Z|Y]|X] = E[Z|X].$$

$$\begin{aligned} \int_{[u(X) \in G]} E[Y|X] dP &= \int_{[X \in u^{-1}(G)]} E[Y|X] dP \\ &= \int_{[X \in u^{-1}(G)]} Y dP \\ &= \int_{[u(X) \in G]} Y dP \\ &= \int_{[u(X) \in G]} E[Y|u(X)] dP \end{aligned}$$

Exercise:  $X$  has countable range...

**Lemma 4.**  $\text{Cov}(u(X), Y - E[X]) = 0$ .

*Proof.*

$$\text{Cov}(u(X), Y - E[Y|X]) = E[u(X) E[Y|X]]$$

□

$$\begin{aligned} E[(Y - f(X))^2] &= E[(Y - E[Y|X] + E[Y|X] - f(X))^2] \\ &= E[(Y - E[Y|X])^2] + 2 \text{Cov}(Y - E[Y|X], E[Y|X] - f(X)) + E[f(X)^2] \end{aligned}$$

**Lemma 5.** *The following are equivalent:*

- (1)  $E[Y|X] = Y$
- (2)  $Y = f(X)$  for some  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- (3)  $\text{Cov}(Y, Z - E[Z|X]) = 0$  for all random variables  $Z$ .

*Proof.*

- (1)  $\Rightarrow$  (2)  $E[Y|X]$  is, by definition, a function of  $X$ .
- (2)  $\Rightarrow$  (3) We have:

$$\begin{aligned} \text{Cov}(f(X), Z - E[Z|X]) &= E[f(X)(Z - E[Z|X])] \\ &= E[f(X)Z] - E[f(X) E[Z|X]] \\ &= E[f(X)Z] - E[E[f(X)Z|X]] \\ &= E[f(X)Z] - E[f(X)Z] \\ &= 0. \end{aligned}$$

- (3)  $\Rightarrow$  (1)

□

$$\begin{aligned} E[u(X)Y] &= E[E[u(X)Y|X]] \\ &= E[u(X) E[Y|X]] \end{aligned}$$

Let  $f(x, y)$  be the empirical density associated to the data set  $(x_1, y_1), \dots, (x_n, y_n)$ :

$$f(x, y) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i) \delta(y - y_i)$$

Suppose that  $(X, Y)$  has joint density  $f(x, y)$ . The marginal densities  $f(x)$  and  $f(y)$  of  $X$  and  $Y$  are

$$f(x) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i) \quad \text{and} \quad f(y) = \frac{1}{n} \sum_{i=1}^n \delta(y - y_i).$$

Let's project  $Y - \mathbb{E} Y$  onto the span of the uncorrelated random variables 1 and  $X - \mathbb{E} X$ . It's easy to show (exercise) that  $\mathbb{E} X = \bar{x}$  and  $\mathbb{E} Y = \bar{y}$ .

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i,$$

Therefore,

$$\begin{aligned} \mathbb{E}[(X - \mathbb{E} X)(Y - \mathbb{E} Y)] &= \iint (y - \bar{y})(x - \bar{x}) f(x, y) dx dy \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \text{cov}(x, y). \end{aligned}$$

Obviously,

$$\mathbb{E}[1(Y - \mathbb{E} Y)] = 0.$$

Therefore, the projection of  $Y - \mathbb{E} Y$  onto the span of 1 and  $X - \mathbb{E} X$  is

$$\frac{\mathbb{E}[1(Y - \mathbb{E} Y)]}{\mathbb{E}[1^2]} 1 + \frac{\mathbb{E}[(X - \mathbb{E} X)(Y - \mathbb{E} Y)]}{\mathbb{E}[(X - \mathbb{E} X)^2]} (X - \mathbb{E} X) = \frac{\text{cov}(x, y)}{\text{var}(x)} (X - \bar{x})$$

It follows that the linear regression of  $Y$  on  $X$  is

$$\hat{Y} = \frac{\text{cov}(x, y)}{\text{var}(x)} (X - \bar{x}) + \bar{y}$$

Consider the probability space

$$(\mathbb{R}^2, f(x, y) dx dy),$$

where  $f(x, y)$  is the *empirical density* associated to the data set  $(x_1, y_1), \dots, (x_n, y_n)$ :

$$f(x, y) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i) \delta(y - y_i).$$

Let

$$V := L^2(\mathbb{R}^2, f(x, y) dx dy) = \left\{ Z : \mathbb{R}^2 \rightarrow \mathbb{R} : \iint |Z(x, y)|^2 f(x, y) dx dy \right\}$$

- You want to “average away” the noise. Interpolating noisy data gives wiggly graphs.
- large oscillations near left and right endpoints
- Increasing size of training set increases model complexity (degree).

#### 4. BIAS-VARIANCE DECOMPOSITION

Let  $\hat{\theta} = \hat{\theta}(X)$  be an estimator of  $\theta$ . The *bias of  $\hat{\theta}$*  is defined by

$$\text{Bias}(\hat{\theta}, \theta) = \mathbb{E} \hat{\theta} - \theta.$$

The variance of the random variable  $\hat{\theta}$  is given, as usual, by

$$\text{Var} \hat{\theta} = \mathbb{E} \left[ (\hat{\theta} - \mathbb{E} \hat{\theta})^2 \right]$$

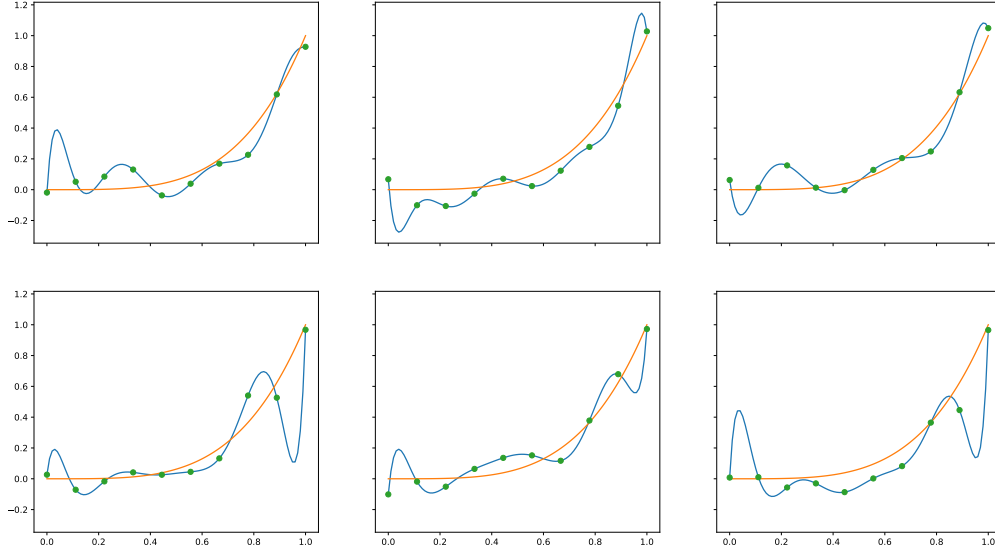


FIGURE 1. —  $y = x^4$ ,  $\bullet$   $y_i = x_i^4 + \text{noise}$ , — polynomial through  $(x_i, y_i)$

**Theorem 6** (Bias-Variance decomposition).

$$\mathbb{E} \left[ (\hat{\theta} - \theta)^2 \right] = \text{Bias}(\hat{\theta}, \theta)^2 + \text{Var} \hat{\theta}$$

*Proof.*

$$\begin{aligned} \mathbb{E} \left[ (\hat{\theta} - \theta)^2 \right] &= \mathbb{E} \left[ (\hat{\theta} - \mathbb{E} \hat{\theta} + \mathbb{E} \hat{\theta} - \theta)^2 \right] \\ &= \mathbb{E} \left[ (\hat{\theta} - \mathbb{E} \hat{\theta})^2 \right] + 2 \mathbb{E} \left[ (\hat{\theta} - \mathbb{E} \hat{\theta})(\mathbb{E} \hat{\theta} - \theta) \right] + \mathbb{E} \left[ (\mathbb{E} \hat{\theta} - \theta)^2 \right] \\ &= \text{Var} \hat{\theta} + \text{Bias}(\hat{\theta}, \theta)^2, \end{aligned}$$

as

$$\mathbb{E} \left[ (\hat{\theta} - \mathbb{E} \hat{\theta})(\mathbb{E} \hat{\theta} - \theta) \right] = (\mathbb{E} \hat{\theta} - \theta) \underbrace{\mathbb{E}[\hat{\theta} - \mathbb{E} \hat{\theta}]}_{=0} = 0. \quad \square$$

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be an unknown function and let  $\hat{f} : \mathbb{R} \rightarrow \mathbb{R}$  be a known approximation to  $f$ . Let  $x_0 \in \mathbb{R}$  and suppose that

$$Y = f(x_0) + \varepsilon, \quad \text{where} \quad \mathbb{E}[\varepsilon] = 0.$$



The *squared prediction error* is

$$\begin{aligned}
(f(x_0) - \hat{f}(x_0))^2 &= \mathbb{E} \left[ (\hat{f}(x_0) - f(x_0))^2 \right] \\
&= \mathbb{E} \left[ (\hat{f}(x_0) - Y - \varepsilon)^2 \right] \\
&= \mathbb{E} \left[ (Y - f + f - \hat{f})^2 \right] \\
&= \mathbb{E} \left[ (Y - f)^2 \right] + 2 \mathbb{E} \left[ (Y - f)(f - \hat{f}) \right] + \mathbb{E} \left[ (f - \hat{f})^2 \right] \\
&= \mathbb{E}[\varepsilon^2] + 2\varepsilon \mathbb{E}[f - \hat{f}] + \text{Bias}(\hat{f}, f)
\end{aligned}$$

Let  $\theta \in \mathbb{R}$ , let  $\varepsilon$  be a random variable with  $\mathbb{E}[\varepsilon] = 0$ , and let

$$Y = \theta + \varepsilon.$$

Let  $\hat{\theta}$  be an estimator of  $\theta$  such that  $\hat{\theta}$  and  $\varepsilon$  are independent.

$$\begin{aligned}
\mathbb{E}[(\hat{\theta} - Y)^2] &= \mathbb{E}[(\hat{\theta} - \theta - \varepsilon)^2] \\
&= \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta} + \mathbb{E}\hat{\theta} - \theta - \varepsilon)^2] \\
&= \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})^2] + \mathbb{E}[(\mathbb{E}\hat{\theta} - \theta)^2] + \mathbb{E}[\varepsilon^2] \\
&\quad + 2 \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)] - 2 \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})\varepsilon] - 2 \mathbb{E}[(\mathbb{E}\hat{\theta} - \theta)\varepsilon] \\
&= \text{Var } \hat{\theta} + \text{Bias}(\hat{\theta}, \theta)^2 + \text{Var } \varepsilon
\end{aligned}$$

We have:

- $\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})^2] = \text{Var } \hat{\theta}$
- $\mathbb{E}\hat{\theta} - \theta$  is a constant, so

$$\begin{aligned}
\mathbb{E}[(\mathbb{E}\hat{\theta} - \theta)^2] &= (\mathbb{E}\hat{\theta} - \theta)^2 = \text{Bias}(\hat{\theta}, \theta)^2, \\
\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)] &= \mathbb{E}[\hat{\theta} - \mathbb{E}\hat{\theta}](\mathbb{E}\hat{\theta} - \theta) = 0 \quad (\text{as } \mathbb{E}[\hat{\theta} - \mathbb{E}\hat{\theta}] = 0), \\
\mathbb{E}[(\mathbb{E}\hat{\theta} - \theta)\varepsilon] &= (\mathbb{E}\hat{\theta} - \theta) \mathbb{E}\varepsilon = 0 \quad (\text{as } \mathbb{E}\varepsilon = 0).
\end{aligned}$$

- $\mathbb{E}[\varepsilon^2] = \text{Var } \varepsilon$
- $\varepsilon$  is independent of  $\hat{\theta}$  and, hence, of  $\hat{\theta} - \mathbb{E}\hat{\theta}$ . Therefore,

$$\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})\varepsilon] = \underbrace{\mathbb{E}[\hat{\theta} - \mathbb{E}\hat{\theta}]}_{=0} \mathbb{E}\varepsilon = 0.$$

If you can sample from a distribution, and you have an unbiased estimator, you can learn the parameters of the distribution. The amount of data you need depends on the variance of the estimator.

## 5. NOTES

Statistics is the science of the *collection*, *analysis*, and *interpretation* of data. [TPE p. 1]

Data analysis: Organization and summarization of data. Emphasize main features. Expose underlying structure. Avoid extraneous assumptions.

Statistical inference: We postulate that the data are values realized by random variables obeying a probability distribution belonging to some known class,  $\mathcal{P}$ . Typically,  $\mathcal{P}$  is indexed by some *parameter space*,  $\Theta$ .

$$\mathcal{P} = \{P_\theta : \theta \in \Theta\}$$

We call the family  $\mathcal{P}$  a *parametric* if  $\Theta \subseteq \mathbb{R}^n$ , for some  $n$ , and *nonparametric*, otherwise. In statistical inference, we use data to infer (point estimation) a plausible value of  $\theta$  or (confidence sets) a subset of  $\Theta$  that plausibly contains  $\theta$ .

The estimation problem: Given  $g : \Theta \rightarrow \mathbb{R}$  and an  $\mathcal{X}$ -valued *random observable*  $X$  distributed according to some  $P \in \mathcal{P}$ , determine  $g(\theta(P))$ . An *estimator* is a function  $\delta : \mathcal{X} \rightarrow \mathbb{R}$ . We want to find an estimator  $\delta$  such that  $\delta(X) \approx g(\theta(P))$ .

A *parametric family of distributions* is one that is naturally indexed by a subset  $\Theta$  of some Euclidean space  $\mathbb{R}^n$ . The set  $\Theta$  is called the *parameter space* of the family.

Suppose we are given a sample space  $\mathcal{X} \subseteq \mathbb{R}^p$  and a family  $\mathcal{P}$  of distributions on  $\mathcal{X}$ .

Let  $X$  be an  $\mathcal{X}$ -valued random vector such that  $X \sim P$  for some unknown  $P \in \mathcal{P}$ .

Using data  $x$  realizing  $X$  to make draw conclusions about  $P$  is called *statistical inference*.

Let  $g$  be a *functional* (real-valued function) on  $\mathcal{P}$ .

Using data  $x$  realizing  $X$  to estimate  $g(P)$  is called *point estimation*.

Point estimation is a type of statistical inference.

Let  $P_{\mu,\sigma}$  be the distribution with density

$$\prod_{i=1}^p \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2}.$$

It's the distribution of an i.i.d. sample of size  $p$  drawn from  $N(\mu, \sigma)$ .

Using such a sample to estimate  $\mu$  (resp.,  $\sigma$ ) is an example of point estimation with

$$\mathcal{X} = \mathbb{R}, \quad \mathcal{P} = \{P_{\mu,\sigma}\} \quad \text{and} \quad g(P_{\mu,\sigma}) = \mu \text{ (resp., } \sigma)$$

A *statistical functional* on  $\mathcal{P}$  is a function  $g : \mathcal{P} \rightarrow \mathbb{R}$ . Let  $\mathcal{P}$  be the set of all probability distributions on  $\mathbb{R}$ . For  $a \in \mathbb{R}$ , define

$$g_a(P) = \int_{-\infty}^a dP(x)$$

$$\mu(P) = \int_{-\infty}^{\infty} x dP(x)$$

$$m_k(P) = \int_{-\infty}^{\infty} (x - \mu(P))^k dP(x)$$

Estimating a functional  $g : \mathcal{P} \rightarrow \mathbb{R}$  from data means constructing a function  $\delta : \mathcal{X} \rightarrow \mathbb{R}$  such that for all distributions  $P \in \mathcal{P}$  and all  $\mathcal{X}$ -valued random variables  $X \sim P$ , the quantity  $\delta(X)$  is “close to”  $g(P)$ . We call  $g$  and  $\delta$  the *estimand* and *estimator*, respectively.

We must make the descriptor “close to” precise if we are to evaluate the quality of an estimator  $\delta$  of a functional  $g$  in any meaningful way. The notion of *bias* is a natural interpretation of closeness. Define

$$\text{Bias}(\delta(X), g(P)) = \mathbb{E} \delta(X) - g(P).$$

If  $\text{Bias}(\delta(X), g(P)) < 0$  (resp.,  $\text{Bias}(\delta(X), g(P)) > 0$ ), then  $\delta(X)$  tends to underestimate (resp., overestimate)  $g(P)$ . We say that  $\delta(X)$  is an *biased* (resp., *unbiased*) *estimator* of  $g(P)$  if  $\text{Bias}(\delta(X), g(P)) \neq 0$  (resp.,  $\text{Bias}(\delta(X), g(P)) = 0$ ).

Let  $X \sim P \in \mathcal{P}$ .

$$\text{Bias}(\delta(X), g(P)) = \mathbb{E}[\delta(X) - g(P)]$$

Note that if  $X \sim P$ , then  $\mathbb{E} \delta(X)$  depends only on  $\delta$  and  $P$  and not on  $X$ :

$$\mathbb{E}[\delta(X)] = \int_{\mathcal{X}} \delta(x) dP(x)$$

Define the (*mean*) *bias* functional associated to the estimator  $\delta$  of  $g$ ,

$$\text{Bias}(\delta, g) : \mathcal{P} \longrightarrow \mathbb{R}$$

by

$$P \mapsto \text{Bias}_P(\delta, g) := \mathbb{E}[\delta(X)] - g(P),$$

where  $X$  is any  $\mathcal{X}$ -valued random variable such that  $X \sim P$ .

$\delta$  is an unbiased estimator of  $g$  if and only if  $\text{Bias}(\delta, g)$  is an unbiased estimator of the zero functional.

Mean bias vs. median bias. Exercise?

## 6. LOGISTIC REGRESSION

Define the *sigmoid function*, also called the *expit function* or the *logistic function*, by

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

It maps  $\mathbb{R}$  bijectively onto  $(0, 1)$ . Its inverse is the *logit function*, defined by

$$\text{logit}(x) = \log \left( \frac{x}{1 - x} \right).$$

The logit function maps  $(0, 1)$  bijectively onto  $\mathbb{R}$ .

$$f(t) = y \log \sigma(t) + (1 - y) \log(1 - \sigma(t))$$

$$\begin{aligned}
f' &= y \frac{\sigma'}{\sigma} + (y-1) \frac{\sigma'}{1-\sigma} \\
&= \frac{1}{\sigma(1-\sigma)} [y\sigma'(1-\sigma) + (y-1)\sigma'\sigma] \\
&= \frac{\sigma'(y-\sigma)}{\sigma(1-\sigma)}
\end{aligned}$$

## 7. NEWTON'S METHOD

7.1. **One variable.** To refine an approximation  $f(a) \approx 0$ , we solve

(first order approximation to  $f(x)$  at  $x = a$ )  $= 0$ .

$$f(a) + f'(a)(x - a) = 0$$

$$\text{refined approximation} = x = a - \frac{f(a)}{f'(a)}$$

Hence, Newton's recursion:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

7.2. **Two variables.** Suppose we want to solve

$$(f(\mathbf{x}), g(\mathbf{x})) = (0, 0) = \mathbf{0}.$$

To refine an approximate solution  $(f(a, b), g(a, b)) \approx (0, 0)$ , we solve the system

(first order approximation to  $f(x, y)$  at  $(x, y) = (a, b)$ )  $= 0$

(first order approximation to  $g(x, y)$  at  $(x, y) = (a, b)$ )  $= 0$ .

$$f(a, b) + f_x(a, b)(x - a) + f_y(a, b)(y - b) = 0$$

$$g(a, b) + g_x(a, b)(x - a) + g_y(a, b)(y - b) = 0$$

$$\begin{pmatrix} f(a, b) \\ g(a, b) \end{pmatrix} + \begin{pmatrix} f_x(a, b) & f_y(a, b) \\ g_x(a, b) & g_y(a, b) \end{pmatrix} \begin{bmatrix} x \\ y \end{bmatrix} - \begin{pmatrix} a \\ b \end{pmatrix} = \mathbf{0}$$

$$\begin{aligned}
\text{refined approximation} &= \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} a \\ b \end{pmatrix} - \begin{pmatrix} f_x(a, b) & f_y(a, b) \\ g_x(a, b) & g_y(a, b) \end{pmatrix}^{-1} \begin{pmatrix} f(a, b) \\ g(a, b) \end{pmatrix} \\
&= \begin{pmatrix} a \\ b \end{pmatrix} - \frac{1}{J_{f,g}(a, b)} \begin{pmatrix} g_y(a, b) & -f_y(a, b) \\ -g_x(a, b) & f_x(a, b) \end{pmatrix} \begin{pmatrix} f(a, b) \\ g(a, b) \end{pmatrix},
\end{aligned}$$

where

$$J_{f,g}(a, b) = \begin{vmatrix} f_x(a, b) & f_y(a, b) \\ g_x(a, b) & g_y(a, b) \end{vmatrix} = f_x(a, b)g_y(a, b) - f_y(a, b)g_x(a, b).$$

Hence, Newton's recursion:

$$\begin{pmatrix} x_{n+1} \\ y_{n+1} \end{pmatrix} = \begin{pmatrix} x_n \\ y_n \end{pmatrix} - \frac{1}{|J_{f,g}(x_n, y_n)|} \begin{pmatrix} g_y(x_n, y_n) & -f_y(x_n, y_n) \\ -g_x(x_n, y_n) & f_x(x_n, y_n) \end{pmatrix} \begin{pmatrix} f(x_n, y_n) \\ g(x_n, y_n) \end{pmatrix}$$

### 7.3. Application to simple logistic regression.

$$\ell(a, b) = \text{log-likelihood}$$

To maximize  $\ell(a, b)$ , we solve

$$\nabla \ell(a, b) = \begin{pmatrix} \frac{\partial \ell}{\partial a} \\ \frac{\partial \ell}{\partial b} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

We apply the above with  $f = \frac{\partial \ell}{\partial a}$  and  $g = \frac{\partial \ell}{\partial b}$ . In this case, the intervening matrix is the *Hessian matrix* of  $\ell$ :

$$H_\ell = \begin{pmatrix} \frac{\partial^2 \ell}{\partial a^2} & \frac{\partial^2 \ell}{\partial a \partial b} \\ \frac{\partial^2 \ell}{\partial b \partial a} & \frac{\partial^2 \ell}{\partial b^2} \end{pmatrix}.$$

Newton's recursion becomes:

$$\begin{pmatrix} a_{n+1} \\ b_{n+1} \end{pmatrix} = \begin{pmatrix} a_n \\ b_n \end{pmatrix} - H_\ell(a_n, b_n)^{-1} \nabla \ell(a_n, b_n)$$

By equality of mixed partials, the *Hessian (determinant)* of  $\ell$  is given by:

$$J_{\nabla \ell} = \frac{\partial^2 \ell}{\partial a^2} \frac{\partial^2 \ell}{\partial b^2} - \frac{\partial^2 \ell}{\partial a \partial b}$$

### 7.4. $n$ -variables.

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \mathbf{f}'(\mathbf{x}_n)^{-1} \mathbf{f}(\mathbf{x}_n), \quad \text{where} \quad \mathbf{f}'(\mathbf{x}) = \left( \frac{\partial f_i}{\partial x_j}(\mathbf{x}) \right)_{i,j=1,\dots,n}.$$

## 8. CONVEXITY

$$(1) \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$\ell(\theta) = \sum_{i=1}^n \ell_i(\theta),$$

where

$$\ell_i(\theta) = y^{(i)} \log \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}))$$

$$\begin{aligned} \frac{\partial \ell_i}{\partial \theta_j} &= y^{(i)} \frac{\sigma'(\boldsymbol{\theta}^T \mathbf{x})}{\sigma(\boldsymbol{\theta}^T \mathbf{x})} x_j^{(i)} - (1 - y^{(i)}) \frac{\sigma'(\boldsymbol{\theta}^T \mathbf{x})}{1 - \sigma(\boldsymbol{\theta}^T \mathbf{x})} x_j^{(i)} \\ &= y^{(i)} (1 - \sigma(\boldsymbol{\theta}^T \mathbf{x})) x_j^{(i)} - (1 - y^{(i)}) \sigma(\boldsymbol{\theta}^T \mathbf{x}) x_j^{(i)} \\ &= (y^{(i)} - \sigma(\boldsymbol{\theta}^T \mathbf{x})) x_j^{(i)} \end{aligned}$$

$$\frac{\partial^2 \ell_i}{\partial \theta_j \partial \theta_k} = -\sigma'(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) x_j^{(i)} x_k^{(i)}$$

$$\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} = -\sum_{i=1}^n \sigma'(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) x_j^{(i)} x_k^{(i)}$$

Let

$$X = \begin{pmatrix} x_j^{(i)} \end{pmatrix} \in \mathbb{R}^{n \times p}, \quad D = \text{diag}(\sigma'(\boldsymbol{\theta}^T \mathbf{x}^{(1)}), \dots, \sigma'(\boldsymbol{\theta}^T \mathbf{x}^{(n)})).$$

Then

$$H_\ell(\boldsymbol{\theta}) = X^T D X.$$

## 9. GRADIENT DESCENT

$$\boldsymbol{\theta}^{(n+1)} = \boldsymbol{\theta}^{(n)} - \alpha \nabla f(\boldsymbol{\theta}^{(n)})^T$$

$$f(\mathbf{x}) \leq f(\mathbf{a}) + \nabla f(\mathbf{a})(\mathbf{x} - \mathbf{a}) + (\mathbf{x} - \mathbf{a})^T \nabla^2 f(\mathbf{a})(\mathbf{x} - \mathbf{a})$$

Descent lemma: Let  $g(t) = f(\mathbf{x} + t\mathbf{y})$

$$\begin{aligned} f(\mathbf{x} + \mathbf{y}) - f(\mathbf{x}) &= g(1) - g(0) \\ &= \int_0^1 g'(t) dt \\ &= \int_0^1 \nabla f(\mathbf{x} + t\mathbf{y}) \mathbf{y} dt \\ &= \int_0^1 (f(\mathbf{x}) + \nabla f(\mathbf{x} + t\mathbf{y}) - f(\mathbf{x})) \mathbf{y} dt \\ &= \int_0^1 \nabla f(\mathbf{x}) \mathbf{y} dt + \int_0^1 (\nabla f(\mathbf{x} + t\mathbf{y}) - \nabla f(\mathbf{x})) \mathbf{y} dt \\ &\leq \nabla f(\mathbf{x}) \mathbf{y} + \int_0^1 \|\nabla f(\mathbf{x} + t\mathbf{y}) - \nabla f(\mathbf{x})\| \|\mathbf{y}\| dt \\ &\leq \nabla f(\mathbf{x}) \mathbf{y} + \int_0^1 L \|t\mathbf{y}\| \|\mathbf{y}\| dt \\ &= \nabla f(\mathbf{x}) \mathbf{y} + \frac{L}{2} \|\mathbf{y}\|^2 \end{aligned}$$

Replace  $\mathbf{y}$  by  $\mathbf{y} - \mathbf{x}$ :

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \frac{1}{L} \nabla f(\mathbf{x}_n)^T$$

$$f(x_{n+1}) \leq f(x_n) - \frac{1}{L} \|\nabla f(x_n)\|^2 + \frac{1}{2L} \|\nabla f(x_n)\|^2 = f(x_n) - \frac{1}{2L} \|\nabla f(x_n)\|^2$$

Being a decreasing sequence that is bounded below,  $f(x_n)$  converges.

$$f(x_0) - f(x_n) = \sum_{k=0}^{n-1} (f(x_k) - f(x_{k+1})) \geq \frac{1}{2L} \sum_{k=0}^{n-1} \|\nabla f(x_k)\|^2$$

$$n \min_{k < n} \|\nabla f(x_k)\|^2 \leq \sum_{k=0}^{n-1} \|\nabla f(x_k)\|^2 \leq 2L(f(x_0) - f(x_n)) \leq 2L(f(x_0) - f(x_\infty))$$

$$\min_{k < n} \|\nabla f(x_k)\|^2 \leq \frac{2L(f(x_0) - f(x_\infty))}{n} = O(1/n)$$

$$\liminf_{n \rightarrow \infty} \|\nabla f(x_k)\|^2 = 0.$$