

SML

Signal Analysis,
Models, and
Machine Learning

A Summary

Part I Linear Signal Representation and Approximation

Chapter 1: Hilbert Spaces

A Hilbert space is a vector space with some additional properties.

1.1 Vector Spaces

Is a generalization of the familiar space \mathbb{R}^n of real n -tuples.

Def 1.1 Let $\mathbb{F} \stackrel{\text{def}}{=} \mathbb{R}$ or $\mathbb{F} \stackrel{\text{def}}{=} \mathbb{C}$. A vector space over \mathbb{F} is a set V with properties:

1. Operation "+" : $V \times V \rightarrow V$
 - a) $(u+v)+w = u+(v+w)$ associative law
 - b) $v+w = w+v$ commutative law
2. There exists $0_V \in V$ s.t. $0_V + v = v \quad \forall v \in V$
3. There exists $(-v) \in V$ s.t. $v + (-v) = 0_V \quad \forall v \in V$
4. Operation "multiplication" $\mathbb{F} \times V \rightarrow V$
 - a) $a(v+w) = av + aw$ distributive law
 - b) $(a+b)v = av + bv$
 - c) $(ab)v = a(bv)$
 - d) $1 \cdot v = v$

Elements of V are called vectors, elements of \mathbb{F} are scalars.

Thm 1.1 For any $v \in V$ and $a \in \mathbb{F}$ the following are true

- | | |
|---|------------------------------------|
| 1. There is only one zero element 0_V in V | 3. $0 \cdot v = a \cdot 0_V = 0_V$ |
| 2. There is only one $v' \in V$ s.t. $v + v' = 0_V$ | 4. $(-a)v = a(-v) = - (av)$ |

A subspace of a vector space V is a subset of V that is a vector space in its own right (with same operations as V and same field \mathbb{F}).

Thm 1.2 Subspace Test A nonempty subset U of some vector space V over \mathbb{F} is a subspace of V iff the following conditions are satisfied:

1. $au \in U \quad \forall a \in \mathbb{F}, u \in U$
2. $u+u' \in U \quad \forall u, u' \in U$

Proof: Take an U and check all props. defined in Def 1.1

Thm 1.3 Intersection Thm for Vector Spaces If U_1, U_2, \dots are subspaces of V , then $U_1 \cap U_2 \cap \dots$ is also a subspace of V .

1.2 Linear Independence and Basis

Let $S = \{v_1, v_2, \dots\} \subset V$ be a countable subset of some vectorspace V over \mathbb{F} .

The subspace spanned by S is the set:

$$\text{span}(v_1, v_2, \dots) \stackrel{\Delta}{=} \left\{ \sum_n a_n v_n : a_n \in \mathbb{F} \right\}$$

(all linear combinations of the vectors in S)

Thm 1.4 $\text{span}(v_1, v_2, \dots)$ is a subspace of V

A vector space V is called finite-dimensional if $V = \text{span}(S)$ for some finite set $S \subset V$. The vectors $v_1, v_2, \dots \in V$ are called linearly independent if

$$a_1 v_1 + a_2 v_2 + \dots = 0_V$$

with $a_i \in \mathbb{F}$ has only the solution $a_1 = a_2 = \dots = 0$. Linear dependent if \exists another solution.

Thm 1.5 Unique Representation Thm Let V be some vector space over \mathbb{F} and let $v_1, \dots \in V$ be lin. indep. If $v \in V$ can be written as

$$v = \sum_n a_n v_n = \sum_k b_k v_k$$

for some $a_n, b_k \in \mathbb{F}$, then $a_k = b_k$ for $k = 1, \dots, n$

A basis of V over \mathbb{F} is a countable set $\{v_i, \dots\} \subset V$ of lin. indep. vectors, s.t. $V = \text{span}(v_1, \dots)$

Eg $\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ is the basis of the vector space \mathbb{F}^3 .

Thm 1.6 Basis Extension Thm Any set of lin. indep. vectors $\{v_1, v_2, \dots, v_m\} \subset V$ can be extended to a basis of V

Thm 1.7 If both v_1, v_2, \dots, v_m and w_1, w_2, \dots, w_n form a basis of V , then $m=n$

$\dim V$ is the dimension of V , the unique number of elements in any basis of V

1.3 Inner Product and Norm notation: complex conjugate: \bar{x} , $A^H \triangleq \bar{A^T}$

Def 1.2 V a vector space over \mathbb{F} . Inner product is a function $V \times V \rightarrow \mathbb{F}: (v, w) \mapsto \langle v, w \rangle$

- | | | |
|---|---|---|
| 1. $\langle u+v, w \rangle = \langle u, w \rangle + \langle v, w \rangle$ | 2. $\langle v, w \rangle = \overline{\langle w, v \rangle}$ | 3. $\langle av, w \rangle = a \langle v, w \rangle$ |
| 4. $\langle v, v \rangle \geq 0$ | 4.a $\langle v, v \rangle = 0 \Leftrightarrow v = 0_V$ | $\ v\ \triangleq \sqrt{\langle v, v \rangle}$ |
| 5. $\langle u, v+w \rangle = \langle u, v \rangle + \langle u, w \rangle$ | 6. $\langle v, aw \rangle = \bar{a} \langle v, w \rangle$ | 7. $\langle v, 0_V \rangle = 0$ |

Some examples: In \mathbb{F}^n $\langle v, w \rangle \triangleq \sum v_k \bar{w}_k = w^H v$ generally $\langle v, w \rangle \triangleq \sum a_k v_k \bar{w}_k = w^H A v$

$$A \triangleq \text{diag}(a_1, \dots, a_n)$$

In $\mathbb{Z} \rightarrow \mathbb{F}$ $\langle f, g \rangle \triangleq \sum f[k] \bar{g}[k]$ $T \rightarrow \mathbb{F}, T = \mathbb{R}$ $\langle f, g \rangle = \int_T f(t) \bar{g}(t) dt$

Thm 1.8 Schwarz Inequality $|\langle v, w \rangle| \leq \|v\| \cdot \|w\|$ for all $v, w \in V$

Thm 1.9 1. $\|v\| \geq 0$ 2. $\|v+w\| \leq \|v\| + \|w\|$ 3. $\|av\| = |a| \cdot \|v\| \rightarrow \|v-w\| \leq \|v\| + \|w\|$
 $\|v\|=0 \Leftrightarrow v=0_V$

1.4 Technicalities

The inner product $\langle f, g \rangle = \sum_{k \in \mathbb{Z}} f[k] \overline{g[k]}$ may fail to converge. Or in $\langle f, g \rangle = \int_T f(t) \overline{g(t)} dt$ if $\|f\| = \int_T |f(t)|^2 dt = 0$ we cannot conclude that $f(x) = 0$ for all $x \in T$. This problem is solved by the following trick: For any $v \in V$, let $[v] = \{v' \in V : \|v - v'\| = 0\}$ is an equivalence class. The vector space V is then replaced by the space $\{[v] : v \in V\}$.

- Have to redefine equality: " $f = g$ " means $\int_T |f(t) - g(t)|^2 dt = 0$
- " $x = y$ " means $E[|x - y|^2] = 0$ equiv. $P(x = y) = 1$
- " $X = Y$ " with probability one (w.p.1)

Thm 1.10 Let $T = \mathbb{R}$ or $T = \mathbb{C}$ and V_{fin} set of all functions $f: T \rightarrow T$ that satisfies $\sum_{k \in \mathbb{Z}} |f[k]|^2 < \infty$

1. V_{fin} is a vector space (a subspace of V)
2. The inner product $\sum_{k \in \mathbb{Z}} f[k] \overline{g[k]}$ is well defined and finite $\forall f, g \in V_{fin}$

Thm 1.11 V : space of all T -valued r.v. $V_{fin} \subset V$ set of those r.v. X for which both $E[X]$ and $E[X^2]$ are well defined and finite

1. V_{fin} is a vector space (subspace of V)
2. Inner product (correlation) $E[X \overline{Y}]$ is well defined and finite $\forall X, Y \in V_{fin}$

1.4.3 Hilbert Spaces

A vector space with an inner product is a Hilbert space iff an additional cond. is met:

Def 1.3 A vector space V with an inner product is a Hilbert space if every Cauchy sequence of elements in V converges to ^{an} element of V .

Cauchy sequence: $v_1, v_2, \dots \in V$. for every $\epsilon > 0$ there exists a positive integer n_0 s.t. $\|v_k - v_m\| < \epsilon$ holds for all $m > n_0, k > n_0$.

Ex 1.11 Set of cont. fun. $\mathbb{R} \rightarrow \mathbb{R}$ with inner product, finite norm is a vector space. But

$$f_n(x) = \begin{cases} 0 & x \leq 0 \\ nx & 0 \leq x \leq \frac{1}{n} \\ 1 & x > \frac{1}{n} \end{cases}$$

for $n = 1, 2, \dots$ form a Cauchy seq. of cont. fun. which converges to the step function , which is not cont. and not in $V \rightarrow$ no Hilbert space.

Some Hilbert Spaces:

- T^n
- The set $\ell_2^2(\mathbb{Z})$ of functions $\mathbb{Z} \rightarrow T$ (discrete time signals) with fin. energy
- $L_T^2(\mathbb{R})$ square integrable (finite energy) Radon $\mathbb{R} \rightarrow T$
- $L_T^2(T)$, T an interval of \mathbb{R}
- Set of T -valued random variables with finite first and second moments (mean & variance)

1.5 Orthogonal Vectors and Orthonormal Bases

Two vectors in a HS (Hilbert Space) are orthogonal if $\langle v, w \rangle = 0$. We then have

$$\|v+w\|^2 = \|v\|^2 + \|w\|^2 \quad \|v-w\|^2 = \|v\|^2 + \|w\|^2 \leftarrow \text{Pythagoras' Theorem}$$

For random variables: $E[|x+y|^2] = E[|x|^2] + E[|y|^2]$

Thm 1.13 Let v_1, v_2, \dots be pairwise orthogonal (ie $\langle v_u, v_m \rangle = 0$ if $u \neq m$) and $v_n \neq 0$ $n=1, 2, \dots$. Then v_1, v_2, \dots are linearly independent.

A subset B of a HS is orthonormal if $\langle v, w \rangle = \begin{cases} 1, & v=w \\ 0, & v \neq w \end{cases}$ for all $v, w \in B$. An orthonormal basis is a basis of orthonormal vectors.

Thm 1.14 Orthonormal Representation Then Let V be a HS with orthonormal basis B . Then

$$v = \sum_{w \in B} \langle v, w \rangle w \quad \text{for every } v \in V$$

Thm 1.15 Parseval's Theorem Let V be a HS with orthonormal basis B : $\langle u, v \rangle = \sum_{w \in B} \langle u, w \rangle \overline{\langle v, w \rangle}$ $\forall u, v \in V$

Thm 1.16 u_1, u_2, \dots orthonormal vecs in HS V , $U = \text{span}(u_1, u_2, \dots)$. $x \in V$ then $v \triangleq x - \sum_k \langle x, u_k \rangle u_k$ is orthogonal to every $u_k \in U$

Thm 1.17 Gram-Schmidt Orthonormalization Let V be a vector space and y_1 lin. indep. in V .

Let u_1 be defined as $u_1 \triangleq \frac{\tilde{u}_1}{\|\tilde{u}_1\|}$ with $\tilde{u}_1 \triangleq y_1$ and $\tilde{u}_k = y_k - \sum_{i=1}^{k-1} \langle y_k, u_i \rangle u_i$ for $k > 1$ then the set $\{u_1, \dots, u_n\}$ is an orthonormal basis of $\text{span}(y_1, \dots, y_n)$ for $n=1, 2, \dots$

Ex (Prob 2.3) Find an orthonormal basis for the subspace U of \mathbb{R}^3 spanned by $y_1 = (1, \sqrt{2}, 1)^T$ $y_2 = (1, 0, 1)^T$

$$u_1 = \frac{\tilde{u}_1}{\|\tilde{u}_1\|} = \frac{y_1}{\|y_1\|} = (1, \sqrt{2}/\sqrt{3}, 1)^T \quad \tilde{u}_2 = y_2 - \langle y_2, u_1 \rangle u_1 = (1, -1/\sqrt{3}, 1)^T \quad u_2 = \frac{\tilde{u}_2}{\|\tilde{u}_2\|} = (1, -1/\sqrt{3}, 1/\sqrt{3})^T$$

1.6 Orthogonal Complement Let $F=\mathbb{R}/F=\mathbb{C}$ and U be a subspace of some HS V over F .

The orthogonal complement U^\perp of U is the set $U^\perp = \{v \in V : \langle v, u \rangle = 0 \text{ for all } u \in U\}$

Thm 1.18 The orthogonal complement U^\perp is a subspace of V and $U \cap U^\perp = \{0\}$.

Thm 1.19 Let $U = \text{span}(y_1, \dots) \subset V$. Then for any $v \in V$ we have $v \in U^\perp \iff \langle v, y_k \rangle = 0, k=1, 2, \dots$

Thm 1.20 Orthogonal decomposition Then Let V be a HS and U a subspace of V spanned by $y_1, \dots, y_r \in V$. Then every $x \in V$ has a unique decomposition into $x = u+v$ $u \in U$ $v \in U^\perp$ implies $\dim U + \dim U^\perp = \dim V$

1.7 Projection to a Subspace Let U a subspace of HS V . The projection of $x \in V$ to U is the unique vector u in the decomposition $x = u + v$, $u \in U$ $v \in U^\perp$.

Thm 1.21 The projection $V \rightarrow U : x \mapsto \hat{x}$ is a linear mapping, ie, $ax \mapsto a\hat{x}$ and $x_1 + x_2 \mapsto \hat{x}_1 + \hat{x}_2$

Thm 1.22 Orthogonality Principle Let $U = \text{span}(y_1, y_2, \dots)$ a subspace of some HS V . Then $u \in U$ if and only if the projection of $x \in V$ to U if and only if $\langle x - u, y_k \rangle = 0, k=1, 2, \dots$

Thm 1.23 Orthonormal-Projection Then V a HS, U subspace of V with orthonormal basis $\{u_1, u_2, \dots\}$. Then the projection of $x \in V$ to U is $u = \sum_k \langle x, u_k \rangle u_k$

Thm 1.24 Orthonormal Projection in Lin. Alg. For $F=\mathbb{R}/F=\mathbb{C}$ let $B = (b_1, b_2, \dots) \in F^{m \times n}$ with orthonormal columns $b_1, \dots, b_n \in F^m$ $b_k^H b_e = \begin{cases} 1 & k=e \\ 0 & k \neq e \end{cases}$ Then the projection of $x \in F^m$ to the space spanned by b_1, b_2, \dots, b_n is $u = BB^H x$

$$\text{Proof: } u = \sum b_k \langle x, b_k \rangle = \underbrace{\sum b_k b_k^H}_{B} x = \underbrace{\begin{pmatrix} b_1 & \dots & b_n \end{pmatrix}}_{B} \begin{pmatrix} b_1^H \\ \vdots \\ b_n^H \end{pmatrix} = BB^H x$$

A Basic concepts of Probability Theory

Probability space (Ω, \mathcal{E}, P) :

- Ω : experimental outcomes
- \mathcal{E} : set of events (=subset of Ω) closed under unions and intersections
- $P: \mathcal{E} \rightarrow [0,1]$ probability measure

Discrete random variable: is a function $X: \Omega \rightarrow S$ for every $s \in S$, $\{\omega \in \Omega : X(\omega) = s\}$ is an event.
By "the event $X=s$ " we mean the set (event) $\{\omega \in \Omega : X(\omega) = s\}$.

Probability mass function: For some discrete RV X with codomain S is the function $S \rightarrow [0,1]: s \mapsto P(X=s)$ $p_X(\cdot) = p(\cdot)$

$$p_X(x) \triangleq P(X=x) \quad p_{X,Y}(x,y) \triangleq P(X=x \text{ and } Y=y) \quad p_{X|A}(x) \triangleq P(X=x | A) \quad p_{X|Y}(x|y) \triangleq P(X=x | Y=y)$$

$$\sum_y p(x,y) = p(x) \quad \text{Independence: if } p(x,y) = p(x) \cdot p(y)$$

$$\text{Conditional probability: } P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(A) = \sum_i P(A|B_i) \cdot P(B_i) \quad P(B_i|A) = \frac{P(A|B_i) \cdot P(B_i)}{P(A)}$$

Distribution Function: $F_X: \mathbb{R} \rightarrow \mathbb{R}$ $F_X(r) = P(X \leq r)$ density: $f_X(x) \triangleq \frac{d}{dx} F_X(x) \quad \int_{-\infty}^r f_X(y) dy = F_X(r)$

$$p(x,y,z) = p(x)p(y|x)p(z|x,y) \quad E[X] = \sum_y E[X|Y=y] \cdot p(y) \quad P(A|Y=y) = \frac{P(A) p_{Y|A}(y)}{p_Y(y)} \quad P(A) = \sum_y p_Y(y) P(A|Y=y)$$

$$\text{Mean (Expectation): } m_X = E[X] = \int_{-\infty}^{\infty} x \cdot f_X(x) dx \quad E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx \quad E[\alpha X + \beta Y] = \alpha E[X] + \beta E[Y]$$

X, Y indep. $\Rightarrow X, Y$ uncorrelated $\Leftrightarrow E[X \cdot Y] = E[X] \cdot E[Y]$

$$\text{Variance: } \text{Var}(X) \triangleq E[(X - m_X)^2] = E[X^2] - m_X^2 \quad \text{if } X \text{ coupler: } \text{Var}(X) \triangleq E[(X - m_X)(\bar{X} - m_X)] = E[X^2] - 2m_X E[X] + m_X^2$$

X, Y uncorrelated: $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$ $\text{if } X, Y \text{ orthogonal } (E[X \cdot Y] = 0)$: E

$$\text{Correlation matrix: } R_X \triangleq E[XX^H] \quad \text{Cov. matrix: } V_X \triangleq E[(X - m_X)(X - m_X)^H] = E[XX^H] - m_X m_X^H$$

$$R_{AX} = A R_X A^H \quad V_{AX} = A V_X A^H$$

Weak law of large numbers: If X_1, X_2, \dots are pairwise uncorrelated and identically distributed with finite mean $E[X_k] = m_X$ and finite variance, then

$$\lim_{n \rightarrow \infty} P\left(|m_X - \frac{1}{n} \sum_{k=1}^n X_k| < \varepsilon\right) = 1 \quad \text{for any positive } \varepsilon.$$

Strong law of large numbers: If X_1, \dots iid with finite mean and $E[|X_k|] < \infty$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n X_k = m_X \quad \text{with probability one.}$$

B Concepts & Terms of estimation theory

ML: Maximum Likelihood Estimation Given: probability law $p(y|x)$ and observation $y=y$
 Estimate: $\hat{x} = \underset{x}{\operatorname{argmax}} p(y|x)$ Given samples y , find a distribution P with parameters x that fits best.
 $\underbrace{\quad}_{\text{likelihood function}}$

MAP: maximum a-posteriori Given: joint probability $p(x,y)$ and observation $y=y$

$$\text{Estimate } \hat{x} = \underset{x}{\operatorname{argmax}} p(x|y) = \underset{x}{\operatorname{argmax}} p(x,y) = \underset{x}{\operatorname{argmax}} p(y|x)p(x)$$

Bayesian Estimation Given: joint probability $p(x,y)$, observation $y=y$ and cost function $u(x,\hat{x})$

$$\text{Estimate: } \hat{x} = \underset{x}{\operatorname{argmin}} E[u(x,x')|Y=y] = \underset{x}{\operatorname{argmin}} \sum_x p(x|y)u(x,x') = \underset{x}{\operatorname{argmin}} \sum_x p(x|y)u(x,x')$$

MMSE: Bayesian Min. Mean squared Error Bayesian Estimation with $u(x,\hat{x}) = \|x - \hat{x}\|^2 \rightarrow \hat{x} = E[x|Y=y]$

Penalized-Likelihood Given: $p(y|x)$, preference function $p(x) \geq 0$, observation $y=y$

$$\text{Estimate: } \hat{x} = \underset{x}{\operatorname{argmax}} p(y|x)p(x) = \underset{x}{\operatorname{argmin}} (-\log p(y|x) - \log p(x))$$

Conjugate Priors and Virtual-Observation Priors

A virtual-observation prior is a preference function $p(x)$ of the form $p(x) = p(y_0|x)$ for some hypothetical initial/additional "observation" $y_0=y_0$. A conjugate prior is a prior $p(x)$ s.t. the posterior $p(x|y) \propto p(y|x)p(x)$ has the same functional form as the likelihood function $p(y|x)$

Non-Bayesian Esti. for Mean, Var, Cov. Matx Given: y_1, \dots, y_L iid rv with values in \mathbb{R}^n and unknown mean μ_y and cov. matrix V_y . Let $y^{(1)}, \dots, y^{(L)}$ be realizations of y_i from which we estimate μ_y and V_y .

$\hat{\mu}_y = \frac{1}{L} \sum_{e=1}^L y^{(e)}$ satisfies $E[\hat{\mu}_y] = \mu_y$ and $\lim_{L \rightarrow \infty} \hat{\mu}_y = \mu_y$ w.p.1

Sample covariance matrix $V_S \triangleq \frac{1}{L} \sum_{e=1}^L (y^{(e)} - \hat{\mu}_y)(y^{(e)} - \hat{\mu}_y)^T = \frac{1}{L} \sum_e y^{(e)}(y^{(e)})^T - \hat{\mu}_y \hat{\mu}_y^T$

satisfies $E[V_S] = E[Y Y^T] - E[\hat{\mu}_y \hat{\mu}_y^T] = \frac{L-1}{L} V_y$ and $\lim_{L \rightarrow \infty} V_S = V_y$ w.p.1

If y_1, \dots, y_L are Gaussian, $\hat{V}_y = V_S$.

For small L gets problematic. Regularizations such as $\hat{V}_y = V_S + \frac{\rho}{L} I$, $\rho \in \mathbb{R}$ $\Leftrightarrow \rho = \text{const}$ or $\text{tr}(V_S)$

1.8 Approximation in a Subspace, Least Squares, LMMSE

Least Squares Given: • vector $x = (x_1, \dots, x_m)^T$ with entries in $\mathbb{F} = \mathbb{R}$ or $\mathbb{F} = \mathbb{C}$
• matrix $A \in \mathbb{F}^{m \times n}$ over \mathbb{F}

Problem: Find $h = (h_1, \dots, h_n)^T$ such that $\|Ah - x\|$ is as small as possible.

LMMSE (Linear Minimum Mean Squared Error)

Problem: Wish to estimate a random variable X from random variables Y_1, Y_2, \dots

$$\hat{X} = \sum_h h_h Y_h \quad \text{s.t.} \quad E[\|\hat{X} - X\|^2] \quad \text{is as small as possible.}$$

1.8.1 FIS problem and solution

Let x, y_1, y_2, \dots be vectors in some HS V . Wish to find $\hat{x} = \text{span}(y_1, \dots)$ s.t. $\|\hat{x} - x\|$ is as small as possible.

Thm 1.25 Let $\{y_1, y_2, \dots\}$ be a countable subset of HS V and $U = \text{span}(y_1, \dots)$. For some $x \in V$, let $x = u + v$ with $u \in U$ and $v \in U^\perp$. Then any $\hat{x} \in U$ satisfies

$$\|\hat{x} - x\| \geq \|x - u\|,$$

with equality iff $\hat{x} = u$. The unique best approximation of $x \in V$ by $\hat{x} \in U$ is the projection of x to U .

Proof: $\|\hat{x} - x\|^2 = \|(\hat{x} - u) - (x - u)\|^2 = \|\hat{x} - u\|^2 + \|x - u\|^2$

By (Thm 1.22) the best approx. is $\langle \hat{x} - x, y_h \rangle = 0 \quad h=1, 2, \dots$

1.8.2 Least Squares

Let y_1, y_2, \dots, y_n be cols. of $A \rightarrow \hat{x} = Ah = \sum_{h=1}^n h_y y_h$ (same problem as in 1.8.1)

Recall the inner product of two col.vecs $a \cdot b = a^\top b = b^\top a$ and thus

$$\langle \hat{x} - x, y_h \rangle = \underbrace{y_h^\top (\hat{x} - x)}_{\hat{x} - x} = 0 \quad \text{for } h=1, \dots, n \quad \text{or equiv. } A^\top (Ah - x) = 0 \rightarrow A^\top Ah = A^\top x$$

1.8.3 LMMSE The inner product $\langle x, y \rangle$ of two r.v. is $E[X\bar{Y}]$ thus

$$\langle \hat{x} - x, y_h \rangle = E\left[\left(\sum_{e=1}^n h_e y_e - x\right) \cdot \bar{y}_h\right] = 0 \Leftrightarrow \sum_{e=1}^n E[y_e \bar{y}_h] h_e = E[x \bar{y}_h] \quad h=1, \dots, n \quad \boxed{E[x \bar{y}_h] = 0}$$

$$\text{in matrix notation: } \begin{bmatrix} E[\bar{Y}_1 Y_1] & E[\bar{Y}_1 Y_2] & \dots & E[\bar{Y}_1 Y_n] \\ \vdots & \ddots & \ddots & \vdots \\ E[\bar{Y}_n Y_1] & E[\bar{Y}_n Y_2] & \dots & E[\bar{Y}_n Y_n] \end{bmatrix} \cdot \begin{bmatrix} h_1 \\ \vdots \\ h_n \end{bmatrix} = \begin{bmatrix} E[x \bar{Y}_1] \\ E[x \bar{Y}_2] \\ \vdots \\ E[x \bar{Y}_n] \end{bmatrix}, \quad E[\bar{Y} Y^\top] h = E[\bar{X} \bar{Y}]$$

Common 'trick' for offset: $\hat{x} = h_0 + \sum_{h=1}^n h_y y_h = \sum_{h=0}^n h_y y_h$ with $y_0 \triangleq 1$

Least squares / LMSE Examples

Ex1.15 x, z indep. real r.v. $E[z]=0$ observe $y=x+z \rightarrow \hat{x}=hy$ find h LMSE

$$\text{Sol: } E[\underbrace{hY_h}_{\hat{x}}] = 0 \text{ for } h=1: E[(hY_h - x)y] = 0 = hE[Y^2] - E[xy] \rightarrow hE[Y^2] = E[xy]$$

$$\left. \begin{aligned} E[Y^2] &= E[(x+z)^2] = \dots = E[x^2] + E[z^2] \\ E[xy] &= E[(x+z)x] = \dots = E[x^2] \end{aligned} \right\} \quad h = \frac{E[x^2]}{E[x^2] + E[z^2]}$$

Prob 2.8 X : r.v. with μ_x, σ_{x^2} . Observe $Y=ax+W$ $a \in \mathbb{R}$ $\mu_W=0$ σ_w^2 X id. W. Compute LMSE with offset $\rightarrow E[(x - (y_h + h_0))^2]$ as small as possible

$$\text{Sol: } E[(\hat{x}-x)\bar{Y}_E] = 0 \quad \text{write } Y_E = (1, Y)^T \quad \hat{x} = h_0 + h_1 Y = (h_0, h_1) Y_E = h Y_E$$

$$\left. \begin{aligned} E[Y_E Y_E^T] h &= E[X Y_E] \\ E[Y_E Y_E^T] &= E \begin{bmatrix} 1 & Y \\ Y & Y^2 \end{bmatrix} = \begin{pmatrix} 1 & a\mu_x \\ a\mu_x & a^2(\mu_x^2 + \sigma_{x^2}) + \sigma_w^2 \end{pmatrix} \\ E[X Y_E] &= E \begin{bmatrix} X \\ XY \end{bmatrix} = \begin{pmatrix} \mu_x \\ a(\mu_x^2 + \sigma_{x^2}) \end{pmatrix} \end{aligned} \right\} \quad \left. \begin{aligned} h_0 &= \frac{\mu_x \sigma_w^2}{a^2 \sigma_{x^2} + \sigma_w^2} \\ h_1 &= \frac{a \sigma_{x^2}}{a^2 \sigma_{x^2} + \sigma_w^2} \end{aligned} \right.$$

Prob 2.9a) Formulate Least Squares problem: Fit a straight line through $(p_1^{(1)}, p_2^{(1)}), \dots, (p_1^{(m)}, p_2^{(m)})$
ie find h_0, h_1 st. $\sum_{k=1}^m (p_2^{(k)} - f(p_1^{(k)}))^2$ with $f(x) = h_0 + h_1 x$ is minimal

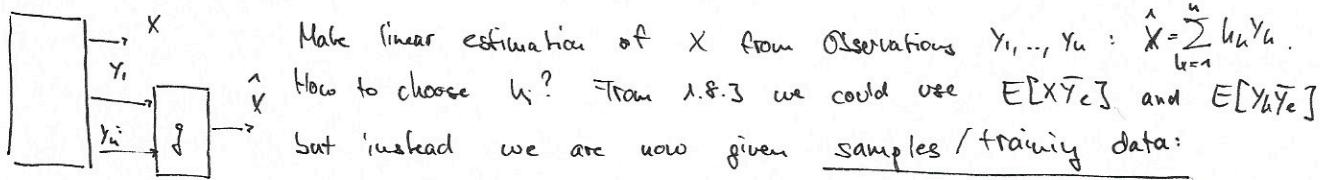
$$\text{Sol: } y^{(k)} \triangleq (1, p_1^{(k)})^T \quad h \triangleq (h_0, h_1)^T \quad f(p_1^{(k)}) = h_0 + h_1 p_1^{(k)} = (y^{(k)})^T h$$

$$x \triangleq (p_2^{(1)}, \dots, p_2^{(m)})^T \quad A \triangleq \begin{pmatrix} y^{(1)T} \\ y^{(2)T} \\ \vdots \\ y^{(m)T} \end{pmatrix} = \begin{pmatrix} 1 & p_1^{(1)} \\ 1 & p_1^{(2)} \\ \vdots & \vdots \\ 1 & p_1^{(m)} \end{pmatrix} \quad \sum_{k=1}^m (p_2^{(k)} - f(p_1^{(k)}))^2 = \|x - Ah\|^2 \text{ which is a LS-problem}$$

Chapter 3: Learning Linear Functions and More About Least-Squares

3.1 Linear-Function Learning (Linear Regression)

3.1.1 Problem and Solution



$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)}) \quad y^{(k)} \triangleq (y_1^{(k)}, \dots, y_n^{(k)})^T$$

We use the vector notation $h \triangleq (h_1, \dots, h_n)^T$ $y \triangleq (y_1, \dots, y_n)^T$ $\boxed{g(y) = \sum_{k=1}^n h_k y_k = h^T y = g^T h}$

Now choose h s.t. the average squared error (ASE) is minimized

$$\text{ASE} = \frac{1}{m} \sum_{k=1}^m \|x^{(k)} - g(y^{(k)})\|^2 = \frac{1}{m} \sum_{k=1}^m \|x^{(k)} - (g^T h)^T h\|^2$$

Writing $x \triangleq (x^{(1)}, \dots, x^{(n)})^T$, $A \triangleq \begin{pmatrix} (y^{(1)})^T \\ (y^{(n)})^T \end{pmatrix}$ this is a least squares problem: $\boxed{\text{w.i. } \text{ASE} = \|x - Ah\|^2} \quad (3.8)$
 $\hookrightarrow A^T Ah = A^T x \quad (3.9)$

3.1.2 Connection with LMMSE

Assume a statistical setting, where the samples $(x^{(k)}, y^{(k)})$ are drawn from a $p(x, y)$ independently.
 A and x defined as in (3.8) then $A^T x = \sum_{e=1}^m y^{(e)} x^{(e)}$ and $A^T A = \sum_{e=1}^m y^{(e)} (y^{(e)})^T$ are random
 $(A^T A h = A^T x)$. By the ^{strong} law of large numbers:

$$\lim_{m \rightarrow \infty} \frac{1}{m} A^T x = E[\sum y^{(e)} x^{(e)}] = E[\bar{Y} x] \quad \lim_{m \rightarrow \infty} \frac{1}{m} A^T A = E[\sum y^{(e)} (y^{(e)})^T] = E[\bar{Y} \bar{Y}^T] \quad \text{w.p.1}$$

It follows that $A^T Ah = A^T x$ converges to $E[\bar{Y} \bar{Y}^T] h = E[\bar{Y} x]$ which coincides with LMMSE.

In a statistical setting, lin-func. learning converges to LMMSE estimation
 in the limit of infinitely many training samples.

C Gaussian Random Variables

$$x \in \mathbb{R}: f_x(z) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(z-\mu_x)^2}{2\sigma^2}} \quad z \in \mathbb{C}: f_z(z) = \frac{1}{\pi\sigma^2} e^{-\frac{|z-\mu_z|^2}{\sigma^2}} \quad z = x + iy: f_z(x+iy) = f_x(x)f_y(y) = \left(\frac{1}{\sqrt{2\pi\sigma_x^2}}\right) e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2}} \left(\frac{1}{\pi\sigma_y^2}\right) e^{-\frac{(y-\mu_y)^2}{2\sigma_y^2}}$$

A general Gaussian random vector $X = (x_1, \dots, x_n)^T$ may be written as $\boxed{X = Ah + u}$
 u : mean vector A : R/C non-singular $n \times n$ matrix $h = (h_1, \dots, h_n)^T$ indep. R/C Gaussian r.v. with $u_i = 0$ $\sigma_u = 1$

$$\text{The covariance matrix } \text{cov}(X) \triangleq V \triangleq A^T A, \quad \boxed{f_X(x) \propto f_U(A^{-1}(x - \mu))} \quad \boxed{f_U(u) \propto \prod_{i=1}^n e^{-\beta/|u_i|^2} = e^{-\beta/|u|^2}} \quad \beta = \begin{cases} 1/2 \text{ real} \\ 1 \text{ complex} \end{cases}$$

$$\text{It follows: } \boxed{f_X(x) \propto e^{-\beta/|A^{-1}(x - \mu)|^2}} = e^{-\beta(A^{-1}(x - \mu))^T A^{-1}(x - \mu)} = \boxed{e^{-\beta(x - \mu)^T (A^T A)^{-1}(x - \mu)}} \quad \begin{cases} V = A^T A \\ V^{-1} = (A^T A)^{-1} \end{cases} \quad \begin{cases} \text{both pos.} \\ \text{definite} \end{cases}$$

$$\text{With } W = V^{-1} = (A^T A)^{-1} \text{ the exponent is in a general quadratic form } q(x) = (x - \mu)^T W(x - \mu) + c \\ = x^T W x - 2x^T W \mu + \mu^T W \mu + c$$

Thm.C.1 Let Y, Z be jointly Gaussian r.v. Then $X \triangleq Y + Z$ is also Gaussian.

Splitting x into two components (y, z) $q(x)$ becomes

$$q(y, z) = ((y - m_y)^H, (z - m_z)^H) \begin{pmatrix} W_{1,1} & W_{1,2} \\ W_{2,1} & W_{2,2} \end{pmatrix} \begin{pmatrix} y - m_y \\ z - m_z \end{pmatrix} + \text{const} , \quad W_{2,1} = W_{1,1}^{-1} \quad (C.17)$$

Thm C.2 (Gaussian Conditioning Thm) If y, z jointly Gaussian with joint dist $\propto e^{-\beta q(y, z)}$, then conditioned on $z=z$ y is Gaussian with mean

$$m_{y|z=z} \stackrel{d}{=} E[Y|z=z] = m_y - W_{1,1}^{-1} W_{1,2} (z - m_z)$$

and covariance matrix $W_{1,1}^{-1}$

Thm C.3 X, Y jointly Gaussian, the MAP estimate of X from $Y=y$ is an affine (=lin. with offset) function of y and coincides both with the MMSE and CMSE (w/ offset) estimate.

Thm C.4 (Gaussian Max/Lut Thm) Let $q(y, z)$ be a quadratic form as in (C.17) with $W_{1,1}$ positive definite. Then $\int_{-\infty}^{\infty} e^{-q(y, z)} dy \propto \max_y e^{-q(y, z)} = e^{-\min_y q(y, z)}$

Thm C.5 (Sum of Quadratic Forms) Let A, B pos. semi-def. matrices. Then

$$\begin{aligned} (x-a)^H A (x-a) + (x-b)^H B (x-b) &= x^H W x - 2 \operatorname{Re}(x^H W_m) + m^H W_m + c \\ W &= A + B \quad m = (A+B)^H (Aa + Bb) \quad c = (a-b)^H A (A+B)^H B (a-b) \\ W_m &= Aa + Bb \end{aligned}$$

3.2 Solving Least-Squares Problems

Problem: Find $h \in F^n$ s.t. $\|Ah-x\|^2$ is minimal. $A \in F^{m \times n}$ $x \in F^m$ (3.17)
 h minimizes iff $A^T Ah = A^T x$

3.2.1 Least Squares by Steepest Descent or gradient descent

Compute gradient of (3.17): $\frac{\partial}{\partial h_i} \|Ah-x\|^2 = \sum_{k=1}^m \frac{\partial}{\partial h_i} (A_{k,i} h_k - x_k)^2 = \dots = 2 \sum_{k=1}^m (A^T)_{k,i} (A_{k,i} h_k - x_k) = 2 (A^T)_{i,:} (Ah - x)$
 $\hookrightarrow \nabla_h \|Ah-x\|^2 \stackrel{d}{=} \left(\frac{\partial}{\partial h_1}, \dots, \frac{\partial}{\partial h_m} \right)^T = 2 A^T (Ah - x)$

Algorithm: $h^{(k+1)} = h^{(k)} + \beta^{(k)} A^T (x - Ah^{(k)})$ with some initial $h^{(0)}$ e.g. $h^{(0)} = 0$ (3.26)
Step size $\beta^{(k)} > 0$

Thm 3.1 For sufficiently small $\beta > 0$ (3.26) with constant $\beta^{(k)} = \beta$ converges to the optimum
 $\hookrightarrow 0 < \beta < \frac{2}{\sigma_{\max}^2}$ σ_{\max}^2 is the largest eigenvalue of $A^T A$.

3.2.2 Stochastic Gradient Descent

If rows of A are sufficiently random, replace (3.26) by

$$h^{(k+1)} = h^{(k)} + \beta^{(k)} (x_k - A_{k,:} h^{(k)}) (A_{k,:})^T$$

\rightarrow in each step use rough estimate of the gradient.

Generally: $h^{(k+1)} = h^{(k)} + \beta^{(k)} \sum_{v \in J_k} (x_v - A_{v,:} h^{(k)}) (A_{v,:})^T \quad J_1 = \{1, \dots, M\}, J_2 = \{M+1, \dots, 2M\} \quad M \ll n$

additionally:

- shuffle rows of A

$$- Smoothing step: h^{(k)} = \frac{1}{L+1} \sum_{j=k-L}^k h^{(j)} \quad L > 0 \text{ window parameter}$$

3.2.3 Speedup by Momentum Term

$$h^{(k+1)} = h^{(k)} + \beta^{(k)} A(x - Ah^{(k)}) + \underbrace{\gamma(h^{(k)} - h^{(k-1)})}_{\text{smoothes trajectory}}$$

Smoothes trajectory and allows larger step size $\beta^{(k)}$ and faster convergence.

3.3 LMSE by Stochastic Gradient Descent (LMS Algorithm)

Statistical setting from 3.1.2 with samples drawn from probability law $p(x, y)$.

We wish to optimize $E[(\hat{x} - x)^2]$ with $\hat{x} = \sum_{e=1}^n h_e y_e$

$$\frac{\partial}{\partial h_e} E[(\hat{x} - x)^2] = \dots = 2 E[(\hat{x} - x) y_e]$$

The update rule becomes

$$h^{(k+1)} = h^{(k)} + \rho^{(k)} E[(x - \hat{x}^{(k)}) (y_1, \dots, y_n)^T]$$

3.4 Overfitting

Problem: We have too little training data for too many parameters. This results in a perfect ASE in training but fails on new data.

i.e. if A has the form $A = [I \ I \ I \ \dots \ I]$ more par. than samples

Example: $y_1 = x + w_1$, $y_2 = w_2$ estimate $\hat{x} = h_1 y_1 + h_2 y_2$ w_1, w_2 zero mean

$$\text{with LMSE we get } h_2 = 0 \quad h_1 = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_{w_2}^2}$$

If we don't have statistics but samples: $(x^{(1)}, y_1^{(1)}, y_2^{(1)}), (x^{(2)}, y_1^{(2)}, y_2^{(2)}) \rightarrow A = \begin{pmatrix} y_1^{(1)} & y_2^{(1)} \\ y_1^{(2)} & y_2^{(2)} \end{pmatrix}, x = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix}$
inverting A and solve for h yields $h_1 = 1 \quad h_2 = -\frac{w_1^{(1)} + w_1^{(2)}}{w_2^{(1)} + w_2^{(2)}} \quad |h_2| \rightarrow 00 \text{ for } w \rightarrow 0$. $\boxed{6}$

Solution: L1 and L2 regularized least squares

3.5 Regularized Least Squares (L2 & L1)

3.5.1 L2-regularized Least Squares

Problem: For given matrix A and col vec x we wish to minimize $\|Ah - x\|^2$

Idea: It often turns out that overfitting is manifested in too large h 's

Solution: Modify least squares problem to minimize $\|Ah - x\|^2 + \lambda^2 \|h\|^2$

Can be reconverted to a standard least squares problem

$$\|Ah - x\|^2 + \lambda^2 \|h\|^2 = \left\| \begin{pmatrix} Ah - x \\ \lambda h \end{pmatrix} \right\|^2 = \left\| \begin{pmatrix} A \\ \lambda \end{pmatrix} h - \begin{pmatrix} x \\ 0 \end{pmatrix} \right\|^2 = \boxed{\|A' h - x\|^2} \quad \boxed{A' = \begin{pmatrix} A \\ \lambda \end{pmatrix}} \quad \boxed{x' = \begin{pmatrix} x \\ 0 \end{pmatrix}}$$

3.5.2 L1-regularized LS (LASSO)

Solution: Replace LS by $\|Ah - x\|^2 + \lambda \sum_{e=1}^n |h_e|$ (3.55) L1 (LASSO): "automatic feature extraction"

- 1. With proper λ , the h that minimizes (3.55) tends to be sparse, i.e. only few nonzero components in h
- 2. (3.55) is convex and quite easily minimized numerically.

Thm 3.2 Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $f(u) \geq 0 \quad \forall u \in \mathbb{R}^n$, diffable with cont. derivatives $\frac{\partial f}{\partial u_e}$ at $u=0$ For $\lambda > 0$ let show to $\hat{h}(\lambda) \triangleq \arg \min_h (f(h) + \lambda \sum_{e=1}^n |h_e|)$ Then for each $e \in \{1, 2, \dots, n\}$ \exists a finite threshold $\lambda_e \in \mathbb{R}$ s.t. $\hat{h}_e(\lambda) = 0$ for $\lambda > \lambda_e$

Stuff:

- Nonzero h_e : corresponding cols of A are support vectors
- First run LASSO to find support vectors, then run $\|Ah - x\|^2$ with only these in A' and x'
- Stochastic Grad. Descent: $h^{(k+1)} = h^{(k)} + \beta (x^{(k)} - \hat{x}^{(k)}) y^{(k)} - \beta \frac{\lambda}{2} \operatorname{sgn}(h^{(k)})$

Chapter 4: SVD and PCA

4.1 SVD: Singular-Value Decomposition

Notation: $z \in \mathbb{C}$: \bar{z} is its complex conjugate

A matrix/vector: A^T is its transpose $A^H = \overline{A^T}$ transpose and complex conjugate

A real square matrix is orthogonal if $AA^T = I$ i.e. $A^{-1} = A^T$

A complex square matrix is unitary if $AA^H = I$ i.e. $A^{-1} = A^H$

If A is (orthogonal or) unitary, then $\|Ax\| = \|x\|$ for all vectors x

$$(\|Ax\|^2 = (Ax)^H(Ax) = x^H A^H A x = x^H x = \|x\|^2)$$

Thm 4.1 (SVD Thm) A \mathbb{R} or \mathbb{C} matrix A can be written as $A = USV^H$ where

• U and V are unitary (real and orthogonal if A is real)

U, V are square
 $U \in \mathbb{R}^{m \times m}$ or $\mathbb{C}^{m \times m}$

• S is diagonal with nonnegative real elements $S = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ & & & 0 \end{bmatrix}$ σ_i : singular values

• condition number $= \frac{\sigma_{\max}}{\sigma_{\min}}$ of a nonsingular square matrix $V^H V = I$ $U^H U = I$

rows \rightarrow # cols

$A \in \mathbb{C}^{m \times n}$

$U \in \mathbb{C}^{m \times m}$

$S \in \mathbb{R}^{m \times n}$

$V \in \mathbb{C}^{n \times n}$

Example: $A = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix}$ SVD: $A = USV^H = \underbrace{\begin{bmatrix} \sqrt{2} & \sqrt{2} \\ \sqrt{2} & -\sqrt{2} \end{bmatrix}}_U \underbrace{\begin{bmatrix} 5 & 0 & 0 \\ 0 & 3 & 0 \end{bmatrix}}_S \underbrace{\begin{bmatrix} \sqrt{2}/\sqrt{18} & 1/\sqrt{18} & 0 \\ 1/\sqrt{18} & -1/\sqrt{18} & 4/\sqrt{18} \\ 2/3 & -2/3 & -1/3 \end{bmatrix}}_{V^H}$

Thm 4.2 (Singular values and eigenvalues) Let $A \in \mathbb{R}^{m \times n} / \mathbb{C}^{m \times n}$ $A = USV^H$ U_e : e -th col of U V_e : e -th col of V

$$AA^H U_e = \sigma_e^2 U_e$$

$$A^H A V_e = \sigma_e^2 V_e$$

cols of U are eigenvectors of AA^H

cols of V are eigenvectors of $A^H A$

$$\text{eigenvalues } \lambda_e = \sigma_e^2$$

4.2 Pseudo-inverse

The Moore-Penrose pseudo-inverse of $A \in \mathbb{R} / \mathbb{C}$ (not necessarily square) with SVD $A = USV^H$ and non-zero singular values σ_e is

$$A^\# \triangleq V S^\# U^H \quad S^\# \triangleq \begin{bmatrix} \text{diag}(\sigma_1^{-1}, \dots, \sigma_e^{-1}) & 0 \\ 0 & 0 \end{bmatrix} \quad S^\#, A^\# \text{ same dim. as } S^H, A^H$$

zero σ_e are not inverted

Ex 4.2 $A = \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_U \underbrace{\begin{pmatrix} 1 \\ 0 \end{pmatrix}}_S \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{V^H}$ $A^\# = \underbrace{(1)}_V \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{S^\#} \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{U^H} = (1 \ 0)$

Usage: - Compute Least Squares "best fit" to a system of lin. eq. that lacks a unique solution

Thm 4.3 $A \in \mathbb{R}^{m \times n} / \mathbb{C}^{m \times n}$

$$\begin{aligned} AA^\# A &= A \\ A^\# A A^\# &= A^\# \\ (A^H)^\# &= (A^\#)^H \end{aligned} \quad \begin{aligned} (AA^H)^\# &= (A^H)^H A^\# \\ A^\# x = 0 &\iff x^H A = 0 \iff A^H x = 0 \end{aligned} \quad (ABC)^\# = C^H B^H A^\#$$

Thm 4.4 (Pseudo-inverse in Closed Form)

$$\text{If } \text{rank}(A) = \#\text{rows} \quad \rightarrow \quad AA^\# = I \quad A^\# = A^H (AA^H)^{-1}$$

$$\text{If } \text{rank}(A) = \#\text{cols} \quad \rightarrow \quad A^H A = I \quad A^\# = (A^H A)^{-1} A^H$$

In general however, there is no closed-form sol. for the pseudo-inverse.

4.3 Projection to the Column Space

If the cols of $A \in \mathbb{R}^{m \times n}/\mathbb{C}^{m \times n}$ are orthogonal, then the projection to the space spanned by them is $\mathbb{F}^m \rightarrow \mathbb{F}^m: x \mapsto AA^H x$

in space ↑ ↑ function / mapping
out space ↓ ↓ input arguments

If not orthogonal, we have Thm 4.5

Thm 4.5 (Projection by SVD) For any $m \times n$ matrix A over $\mathbb{F} = \mathbb{R}/\mathbb{C}$ the linear mapping

$$(4.44) \quad \mathbb{F}^m \rightarrow \mathbb{F}^m: x \mapsto AA^H x \quad A = USV^H \quad A^H \stackrel{\Delta}{=} VS^H U^H$$

is the projection of x to the subspace spanned by the (not necessarily orthogonal) cols of A . If $A = USV^H$ (4.44) can also be written as $\mathbb{F}^m \rightarrow \mathbb{F}^m: x \mapsto U_{\neq} U^H x$ where U_{\neq} consists of the cols of U that correspond to the nonzero singular values.

Thm 4.6 (SVD Projection Energy) Let y_1, \dots, y_n cols of A over $\mathbb{F} = \mathbb{R}/\mathbb{C}$, $A \in \mathbb{F}^{m \times n}$, $A = USV^H$, σ_i sig. values

$$\text{for all cols } u_k \text{ of } U: \sum_{l=1}^n |\langle y_l, u_k \rangle|^2 = \sum_{l=1}^n |u_k^H y_l|^2 = \sigma_k^2$$

Each squared sing. val. σ_k^2 is the total energy of the projections of y_1, \dots, y_n to the 1D space spanned by u_k .

4.4 SVD and Least Squares

(4.52)

(4.53)

4.4.1 Pseudo-inverse sol. of LS Recallly LS-problem: $\hat{u} = \underset{u}{\operatorname{argmin}} \|Au - x\|^2$ iff $A^H A \hat{u} = A^H x \rightarrow \hat{u} = (A^H A)^{-1} A^H x$

Thm 4.7 (Pseudo-inv. sol. of LS) The vector $\hat{u} = A^H x$ is a sol. of (4.53) i.e. minimizes $\|Au - x\|^2$.

Moreover, if u' also minimizes $\|Au - x\|^2$, then $\|u'\| \leq \|u\|$ with equality only if (iff) $u' = u$.

4.4.2 Singular Values in Least Squares

From Thm 7 with $A = USV^H = U \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \ddots \end{bmatrix}$ is $u = A^H x = V \underbrace{\begin{bmatrix} \sigma_1^{-1} & 0 & 0 \\ 0 & \sigma_2^{-1} & 0 \\ 0 & 0 & \ddots \end{bmatrix}}_{(4.52)} U^H x$

↪ Small singular values σ_i have large impact on u . This typically happens with random square matrices.

Overfitting is manifested in such large values of u . One way to alleviate this is setting $\sigma_i^{-1} = 0$ for too small σ_i .

4.5 Principle-Components Analysis (PCA) $\mathbb{F} = \mathbb{R}/\mathbb{C}$

Given: Some data in col. vectors $y_1, \dots, y_N \in \mathbb{F}^M$, mostly stripped off its mean s.t. $\sum_{k=1}^N y_k = 0$

Problem: For some K , $1 \leq K \leq M$ wish to find K -dimensional subspace of \mathbb{F}^M that contains as much of the "energy" of data as possible. \hat{y}_n is proj. of y_n , $\sum_{n=1}^N \|\hat{y}_n\|^2$ as large as possible or equiv. $\sum_{n=1}^N \|y_n\|^2 - \sum_{n=1}^N \|\hat{y}_n\|^2$ as small as possible.

Ansatz: Let B be a matr. over \mathbb{F} whose cols form an orthonormal basis of the desired subspace. $B \in \mathbb{F}^{M \times K}$ and $B^H B = I_K$ (orthonormality). The projection to $\text{span}(\text{cols of } B)$: $\mathbb{F}^M \rightarrow \mathbb{F}^M: y \mapsto \hat{y} = B B^H y$

Total energy of projection: $\sum_{n=1}^N \|\hat{y}_n\|^2 = \sum_{n=1}^N \|B B^H y_n\|^2 = \sum_{n=1}^N \|B^H y_n\|^2$ (4.74), (4.75)

Ausatz (cont.): Now let $A \triangleq (y_1, \dots, y_N) \in \mathbb{F}^{M \times N}$ matrix with "data" and $A = USV^H$ its SVD.

Assume $\sigma_1 > \sigma_2 > \dots > 0$

Theorem 4.9 (SVD-PCA Theorem) The $M \times k$ mtx $B \triangleq U \begin{pmatrix} I_k \\ 0 \end{pmatrix}$ (the first k cols of U) maximizes (4.74) among all $M \times k$ matrices with orthonormal columns.

Chapter 5: Some Basics

5.1 Parametrized Estimators vs. Global statistical Models

What is given	What we can apply	Goal:
Full statistical model $p(x,y)$	MAP ML MMSE	Approximat x from y_1, \dots, y_m
Partial model $p(y x)$	ML some penalized likelihood	$p(x y) = \frac{p(x,y)}{p(y)} \propto p(y x)p(x)$

If no statistical information is given, we need samples $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

Naive approach: \mathcal{X}, \mathcal{Y} sets of possible values of x, y . If \mathcal{X}, \mathcal{Y} finite with not too many elements, approximate $p(x,y)$ by $p_s(x,y)$ the empirical distribution of training samples

$$p_s(x,y) = \frac{1}{m} \times (\text{\#occurrences of } (x,y) \text{ in training samples}) \quad \lim_{m \rightarrow \infty} p_s(x,y) = p(x,y) \text{ w.p.1}$$

Requires that each pair appears sufficiently often.

Approach I: Parametrized Global Statistical Models: = "generative model" parameters

Assume the unknown $p(x,y)$ belongs to a family of prob. laws $p(x,y|\theta)$. Based on training data, form an estimate $\hat{\theta}$ of θ . Then use $p(x,y|\hat{\theta})$ to estimate x .

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(x^{(1)}, y^{(1)}, \dots, x^{(m)}, y^{(m)} | \theta)$$

Treated in Part III

Approach II: Parametrized Estimators

The estimation function $g: \mathcal{Y} \mapsto \hat{x}$ belongs to a family of functions $\{g_\theta\}$ with parameter θ . Based on training data, estimate $\hat{\theta}$ which determines the estimator.

$$\text{Example: } g: (y_1, \dots, y_m) \mapsto \sum_{k=1}^n u_k y_k$$

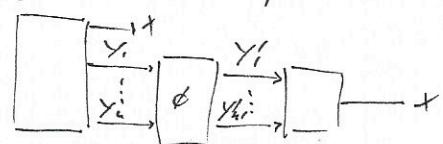
Neural networks \rightarrow Chapter 6 Kernel methods \rightarrow Chapter 7

Nonlinear estimators \rightarrow following

5.2 Features, Transforms, and Lasso Linear Estimation

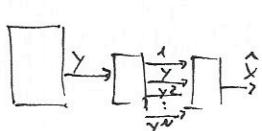
Before estimating, transform observations $Y = (y_1, \dots, y_m)^T$ into $\phi(Y) = Y' = (y'_1, \dots, y'_m)^T$. Dimensions need not be equal

- Benefits:
- Extract features that simplify estimation
 - Reduce dimensionality to a manageable size
 - Transform to higher dimension in which esti is easier
 - Transform nonnumerical Y to numerical Y'



Example 5.1 Polynomial fitting

$$\phi: \mathbb{R} \rightarrow \mathbb{R}^{n+1}: y \mapsto (1, y, y^2, \dots, y^n)^T \quad \hat{x} = \phi(y)^T h$$



Is a Least Squares Problem

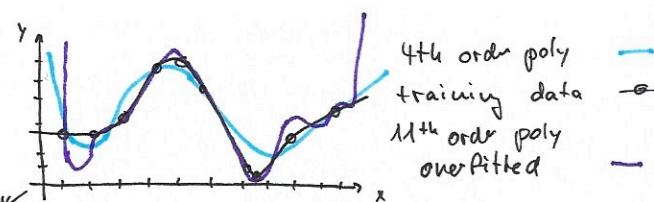
$$\text{ASE} \triangleq \frac{1}{m} \sum_{k=1}^m \|x^{(k)} - \phi(y^{(k)})h\|^2 = \frac{1}{m} \|x - \phi(y)h\|^2$$

$\phi: \mathbb{R} \rightarrow \mathbb{R}^{n+1}: y \mapsto (1, y, y^2, \dots, y^n)^T$

$$\hat{x} = \phi(y)^T h$$

5.3 On Overfitting and Regularization

Problem: Complex models can be adjusted to fit the training data very well but fail completely on new data.



- General Solution:
1. Split data into two parts
 2. Train models on different complexity on first half
 3. Compare performance of trained models on second half of data. Select the best performing model

- Leave-one-out:
1. Train model on all data except 1 sample
 2. Do this for each sample as the excluded one
 3. Test error on excluded sample
 4. Final score is the average error over all repetitions

5.4 Learning Classifiers and Class Probabilities (Logistic Regression)

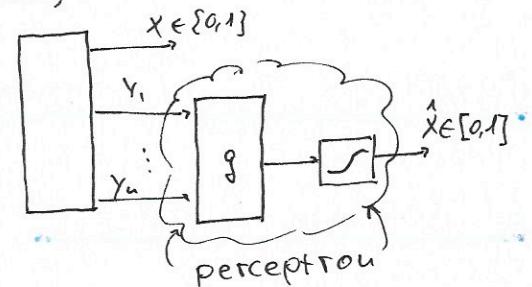
Assumption: X is a $\{0,1\}$ -valued variable. Eg. an indicator of some event or class or condition of interest

Idea: Rather than producing a $\{0,1\}$ estimate of X , produce a $[0,1]$ -valued \hat{x} that may be interpreted as estimate/guess of

Soft estimate/
soft classifier

$$\boxed{P(X=1 | Y=y) = E[X | Y=y]} \quad (5.17)$$

↑
only bcs $X \in \{0,1\}$



5.4.1 Turning Generic Parametrized Estimators into Soft Classifiers

Given: A generic parametrized mapping $g: \mathbb{R}^n \rightarrow \mathbb{R}$ Required: Soft classifier s.t. $\hat{X} \in [0,1]$

Idea: Append g with a function that limits to $[0,1]$ with monotonous increase

E.g. $\boxed{\text{Squa}(x) = \frac{1}{1+e^{-x}}}$

↑
Sgn(x)
0.5
x

Logistic sigmoid function
mapping of this form is called a
perceptron

Example: g might be an affine function \rightarrow soft estimate is $\text{Sgn}(b_0 + b_1 y_1 + b_2 y_2 + \dots + b_n y_n)$

Note For any prob. dist $p(x,y)$ with $\{0,1\}$ -valued X

$$\text{with } L_x(y) \triangleq \ln \frac{P(X=1 | Y=y)}{P(X=0 | Y=y)} = \ln \frac{P(Y=y | X=1)}{P(Y=y | X=0)} + \ln \frac{P(X=1)}{P(X=0)}$$

$$\boxed{P(X=1 | Y=y) = \frac{1}{1 + e^{-L_x(y)}}}$$

5.4.2 Learning by minimizing the ASE

Learning by ASE minimizing works also in this case:

$$E[\text{ASE}] = E\left[\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \hat{x}(y^{(i)})\|^2\right] = E[\|x - \hat{x}(y)\|^2] \quad \lim_{m \rightarrow \infty} \text{ASE} \stackrel{\text{up!}}{=} E[\text{ASE}]$$

(5.25)

(5.25) is minimized by $\hat{x}(y) = E[x | Y=y]$ which is (5.17)

- Probabilities close to one or zero are learned with poor precision
- Minimizing ASE is not a (linear) least-squares problem, even if g is linear

5.4.3 Discriminative Statistical Modeling and ML Learning

Approach: Use statistical models, where samples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ are drawn indep. from a dist. $p(x, y|\theta)$ s.t. $p(x, y|\theta) = p(x|y, \theta)p(y)$ ($p(y)$ assumed not dep. on θ)

Further we define $p(x|y, \theta) \stackrel{def}{=} \begin{cases} \text{sgm}(g(y)) & x=1 \\ 1-\text{sgm}(g(y)) & x=0 \end{cases}$ g a fn. with parameters θ

ML estimation: Now we fit this statistical model to the given samples

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \prod_{e=1}^m p(x^{(e)}, y^{(e)} | \theta) = \dots = \underset{\theta}{\operatorname{argmax}} \frac{1}{m} \sum_{e=1}^m \underbrace{\left(x^{(e)} \ln(\text{sgm}(g(y^{(e)}))) + (1-x^{(e)}) \ln(1-\text{sgm}(g(y^{(e)}))) \right)}_{-K(\theta)} \quad (5.33)$$

$-K(\theta) \leftarrow \text{cost function}$

Minimizing $K(\theta)$ can be done by (stochastic) gradient descent or other methods.

Special case: $g(y) = h_0 + h_1 y_1 + \dots + h_n y_n$ $\theta = h \stackrel{def}{=} (h_0, \dots, h_n)^T$ The gradient ∇_{θ} of K with respect to $\theta=h$

$$\nabla_{\theta} (K(h)) = \frac{-1}{m} \sum_{e=1}^m \left(x^{(e)} (1 - \text{sgm}(h^T y^{(e)})) + (1 - x^{(e)}) \text{sgm}(h^T y^{(e)}) \right)$$

Chapter 6: Neural Networks

Most successful type: Variations of multilayer perceptrons.

6.1 Multilayer Perceptrons

6.1.1 Structure

Two-layer perceptron: $\mathbb{R}^n \rightarrow \mathbb{R} : (y_1, \dots, y_n) \mapsto \xi(y_1, \dots, y_n)$

$$\xi(y_1, \dots, y_n) = g_{\text{out}} \left(w_0 + \sum_{j=1}^m w_j \tilde{g} \left(w_{j,0} + \sum_{i=1}^n w_{j,i} y_i \right) \right)$$

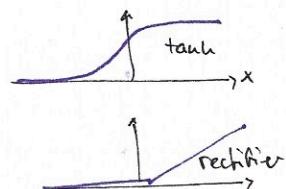
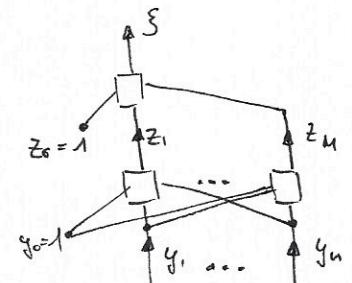
$w_j, w_{j,i}$: real parameters / "weights"

$$y_0 \stackrel{def}{=} 1 \quad z_j \stackrel{def}{=} \tilde{g} \left(\sum_{i=0}^n w_{j,i} y_i \right) \quad z_0 \stackrel{def}{=} 1 \rightarrow \xi = g_{\text{out}} \left(\sum_{j=0}^m w_j z_j \right)$$

$$\tilde{g}(x) = \text{sgm}(x) \quad \text{or} \quad \tilde{g}(x) = \text{tanh}(x) \stackrel{def}{=} \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad \text{or} \quad \tilde{g}(x) = \text{rectifier}(x)$$

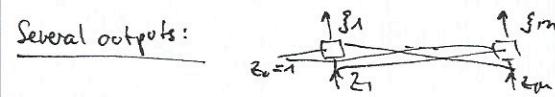
$$g_{\text{out}}(x) = x \quad \text{or} \quad g_{\text{out}}(x) = \text{sgm}(x) \quad \text{if used for classification}$$

↑ rectifier is not at zero



Important All nonlinearities are continuous and (almost) everywhere differentiable, which matters when it comes to optimizing the parameters w_j and $w_{j,i}$.

Hidden nodes: z_1, \dots, z_m are called hidden nodes/variables



Deep networks: With many layers. Most notable class is the convolutional neural networks (CNN)

6.1.2 Learning (training) and Regularization

[in number of samples]

Training means the calculation of weights in all layers from samples $(x^{(1)}, y_1^{(1)}, y_2^{(1)}, \dots, y_n^{(1)})$, $(x^{(2)}, y_1^{(2)}, y_2^{(2)}, \dots, y_n^{(2)}) \dots$

Approach: (Approx.) minimize a cost function such as the $\text{ASE} = \frac{1}{m} \sum_{e=1}^m \|x^{(e)}\|_2^2 (g(y^{(e)}, \theta) - y^{(e)})^2$ on the training samples.

For class probabilities, $K(x)$ from (5.33) can be used as metric

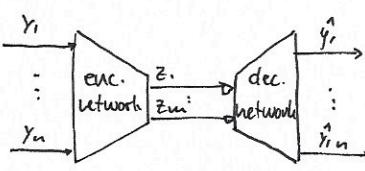
If $g_{\text{out}}(x)=x$ then optimizing w_j is a least squares problem. For the lower layers however, it is not.

Computing layer weights: By back-propagation (not treated). Take the derivative of the cost function with respect to the individual weights and apply a gradient descend.
Can be computationally demanding

- Avoiding Overfitting:
- Switch off random weights for short amount of time
 - Augment training data by distortion/noise
 - Early stopping if error on test data stops decreasing

6.3 Autoencoders, GANs, and Neural-Network Priors

6.3.1 Autoencoder Networks

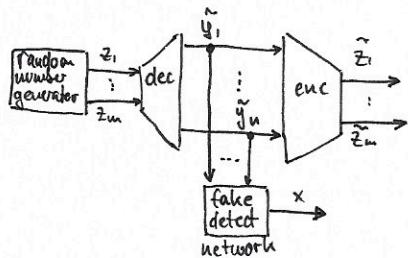


- Is trained to reproduce its input $y \in \mathbb{R}^n$ at the output $\hat{y} \in \mathbb{R}^n$
- $z = (z_1, \dots, z_m)^\top \in \mathbb{R}^m$, where $m < n$ are the hidden units
- enc/dec are some sort of multilayer perceptrons (enc. mostly CNN)
- Trained by unsupervised learning (= on unlabeled data)

- Usages:
- Dim. reduction from \mathbb{R}^n to \mathbb{R}^m , only enc. is used after training
 - Compression: z is used for transmission/storage
 - Denoising: Noise channel on input, noisefree at output \rightarrow use enc. after training
 - Learning a generative model \rightarrow use dec. after learning

Often m is larger than n initially, then use sparsifying cost fun (TSE L1 for example) such that z becomes sparse with $m < n$.

6.3.2 Generative Adversarial Networks (GANs)



- Produces new samples $\hat{y}^{(1)}, \hat{y}^{(2)}, \dots$ according to complex learned statistics
- A way of training autoencoders that can improve their performance
- dec/enc are first trained as an autoencoder
- After this:
 - dec is fed with iid samples $z \in \mathbb{R}^m$ and produce $\hat{y} \in \mathbb{R}^n$
 - fake detect is trained to produce $x=0$ for given samples from enc/dec training y and $x=1$ for output of dec \hat{y}
 - dec trained to produce $x=0$ and $\hat{z}=z$
 - enc trained to produce $\hat{z}=z$

The fake detect network tries to discriminate btw samples of given set y and produced \hat{y} by the dec

The dec tries to fool the fake detect with \hat{y} that cannot be distinguished from y . This arms race results in each of them becoming very good. The enc prevents the dec from producing $\hat{y}=y$.

6.3.3 Neural Network Priors

$$\hat{x} = \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} p(y|x)$$

- The dec. of GANs and autoencoders can be used as priors or regularizers for statistical estimation problems.

$p(y|x)$: likelihood function

- Constraint: x can be produced by decoder network: $f: \mathbb{R}^m \rightarrow \mathbb{R}^n \quad m < n$

$$\hat{x} = f(\hat{z}) \quad \hat{z} = \underset{z}{\operatorname{argmax}} p(y|f(z))$$

6.4 Radial-Basis Function Networks

Is a mapping

$$\mathbb{R}^n \rightarrow \mathbb{R}: y \mapsto \sum_{j=1}^m w_j \phi_j(y) + w_0$$

$w_j \in \mathbb{R}$ parameters/weights
 $\phi_j: \mathbb{R}^n \rightarrow \mathbb{R}$ nonlinear function

Gaussian with mean vect μ_j
 (and corr. matr. W_j^{-1})

$$\phi_j(y) = \exp\left(-\frac{\|y - \mu_j\|^2}{2\sigma^2}\right)$$

$$\phi_j(y) = \exp\left(-\frac{1}{2}(y - \mu_j)^T W_j^{-1} (y - \mu_j)\right)$$

$\phi_j(y) = \text{some function of } \|y - \mu_j\|, y \in \mathbb{R}^n$ parameter/center.

• Training is a nonlin. opt. problem that can be approximated but no guarantee of finding global optimum

• Training is quite different from that of multilayer perceptrons. ϕ_j are trained by samples of y alone (i.e. unlabeled data), optimizing w_j is often a standard least squares problem

K-means clustering algorithm with training samples $y^{(e)}$:

1. Select $K=1$ mean vector $\mu_j, j=1..M$ at random from training set $\{y^{(e)}\}$

2. Cluster the samples into classes C_j by assigning each to the nearest mean vector

$$C_j := \{y^{(e)} : \|y^{(e)} - \mu_j\| \leq \|y^{(e)} - \mu_k\|, k \neq j\}$$

works because

$$\sum_{i=1}^m \|y^{(e)} - \mu_i\|^2$$

decreases in every iteration

→ also must terminate

3. New mean vectors are mean of their clusters: $\mu_j := \frac{1}{|C_j|} \sum_{y \in C_j} y$

4. Repeat 2 & 3 until no change in sets C_j

Chapter 7: Kernel Methods can be used for regression, classification, clustering and more

7.1 Form of Mapping

A kernel based regression function $S \rightarrow \mathbb{R}: y \mapsto \hat{x}$ has the form

$K: S \times S \rightarrow \mathbb{R}$ is called the kernel. The domain S need not be in \mathbb{R}^n (e.g., images, strings, ...)

Example $K(y, y') = \exp\left(-\frac{\|y - y'\|^2}{2\sigma^2}\right)$ Difference with radial-basis function:

• K has no parameters to be adjusted

• Sum (7.1) runs over all training samples

Choosing a suitable K is the most important step in practical applications

7.2 Learning and Regularization

Training data $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

$$\hat{x} = \sum_{j=1}^m w_j K(y^{(e)}, y^{(j)}) + w_0$$

→ LS-problem

#coeff $w_j = \# \text{training samples} \Rightarrow$ can have $\text{ASE}=0$ by

$$x^{(e)} = \sum_{j=1}^m w_j K(y^{(e)}, y^{(j)}) + w_0 \quad l=1..m \text{ for weights}$$

BUT: results in large weights & bad performance on new data. → Regularization needed

↓ L1 regularization

$$\text{ASE}_{L1} = \frac{1}{m} \sum_{e=1}^m \left(x^{(e)} - \sum_{j=1}^m w_j K(y^{(e)}, y^{(j)}) - w_0 \right)^2 + \lambda \sum_{j=1}^m |w_j|$$

can result in zero weights → complexity reduction in (7.2)
 which is essential in kernel methods.

Nonzero $w_j \rightarrow y^{(j)}$ are called support vectors.
 As of section 3.5.2., (7.5) is convex in w_j

7.3 Feature Space Interpretation and Kernel Trick

A nonlinear mapping combined with linear estimation

$$\hat{x} = \phi(y)^T w$$

from section 5.2 for coeff $w \in \mathbb{R}^n$.

If $K(y, z) = \phi(y)^T \phi(z)$ (7.7), (7.1) are essentially equivalent. (7.7) is a Mercer kernel.

Theorem 7.1 (Kernel Trick Thm): $\phi: S \rightarrow \mathbb{R}^n, y^{(1)}, \dots, y^{(m)} \in S$ given and fixed, $K(y, z) = \phi(y)^T \phi(z)$. Then for any given $w_1, \dots, w_m \in \mathbb{R}^n$ there exists $u \in \mathbb{R}^n$ s.t.

$$\sum_{j=1}^m w_j K(y, y^{(j)}) = \phi(y)^T u \quad (7.8) \text{ for all } y \in S.$$

Conversely, for any $u \in \mathbb{R}^n, \exists u' \in \mathbb{R}^n$ and $w \in \mathbb{R}^n$ s.t. $\sum_{j=1}^m w_j K(y, y^{(j)}) = \phi(y)^T u'$ and $\phi(y)^T u' = \phi(y)^T u$ if $\phi(y)$ is in the subspace of \mathbb{R}^n spanned by $\phi(y^{(1)}), \dots, \phi(y^{(m)})$.

Summary: Mappings of form $\hat{x} = \sum w_j K(y, y^{(j)}) + w_0$ are mappings of form $\phi(y)^T w$ in disguise ("rotated")

= Computation and their complexity are quite different

- Sparsifying by L1 doesn't yield the same mappings

Chapter 8: Factor Graphs and Hidden Markov Models

- Statistical models with many variables have usually some modular structure
- These can be represented by factor graphs
- They allow a unified approach to a wide variety of structured models and algorithms

An example modular structure are hidden Markov models.

8.1 Markov Chains and Hidden Markov Models

Def: Three r.v. X, Y, Z form a Markov Chain if, conditioned on event $Y=y$ X and Z are independent \Leftrightarrow depends only on y .

$$\text{ie.: if } p(x, z | y) = p(x|y)p(z|y) \Rightarrow p(x, y, z) = p(x)p(y|x)p(z|y)$$

$$(\text{Generally: } p(x, y, z) = p(x)p(y|x)p(z|x, y))$$

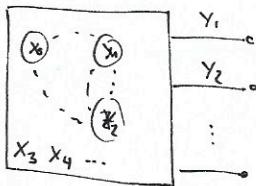
- Holds for prob. mass and prob. density functions

In general: x_1, x_2, \dots, x_n form a MC if, for $1 \leq k \leq n$ conditioned on $x_k = x_k$ (x_1, \dots, x_{k-1}) and (x_{k+1}, \dots, x_n) are independent. i.e. if

$$p(x_1, \dots, x_{k-1}, x_k, \dots, x_n | x_k) = p(x_1, \dots, x_{k-1} | x_k) \cdot p(x_{k+1}, \dots, x_n | x_k)$$

$$\Rightarrow p(x_1, \dots, x_n) = p(x_1)p(x_2|x_1)p(x_3|x_2) \dots p(x_n|x_{n-1})$$

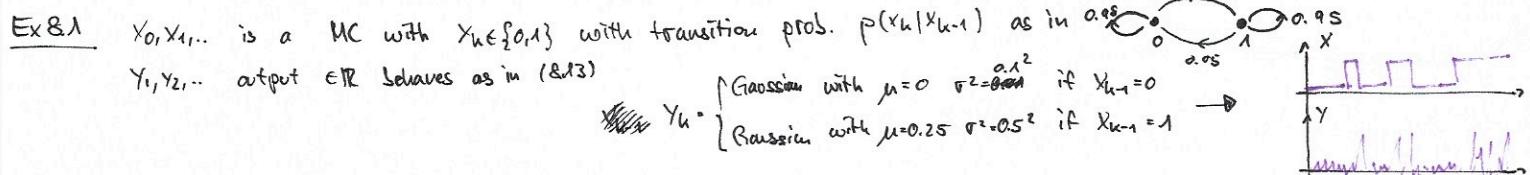
Def (Hidden Markov Model): x_0, x_1, x_2, \dots hidden variables / state variables not directly observable
 y_1, y_2, y_3, \dots observable variables



for all $n \geq 1$ x_0, x_1, \dots, x_n form a Markov Chain and $p(y_1, \dots, y_n | x_0, \dots, x_n) = \prod_{k=1}^n p(y_k | x_k, x_{k-1})$ (8.13)

$$\text{It follows } p(x_0, \dots, x_n, y_1, \dots, y_n) = p(x_0) \prod_{k=1}^n p(x_k, y_k | x_{k-1})$$

It is usually assumed that the r.v. x_i take values in some finite set



Common computational problems For fixed observations $y_1 = y_1, \dots, y_n = y_n$

$$\text{Current state estimation: } p(x_n | y_1, \dots, y_n) = \frac{p(x_n, y_1, \dots, y_n)}{p(y_1, \dots, y_n)} \propto p(x_n, y_1, \dots, y_n) = \sum_{x_0} \dots \sum_{x_{n-1}} p(x_0, \dots, x_n, y_1, \dots, y_n) \quad (8.17)$$

$$\text{Prediction of next output: } p(y_{n+1} | y_1, \dots, y_n) = \frac{p(y_{n+1}, y_1, \dots, y_n)}{p(y_1, \dots, y_n)} \propto p(y_{n+1}, y_1, \dots, y_n) = \sum_{x_0, \dots, x_n} p(x_0, \dots, x_n, y_1, \dots, y_{n+1}) \quad (8.20)$$

$$\text{Estimation of state at any time } k \leq n: p(x_k | y_1, \dots, y_n) = \frac{p(x_k, y_1, \dots, y_n)}{p(y_1, \dots, y_n)} \propto p(x_k, y_1, \dots, y_n) = \sum_{x_0, \dots, x_{k-1}} p(x_0, \dots, x_k, y_1, \dots, y_n) \quad (8.22)$$

$$\text{MAP estimation of state trajectory: } (\hat{x}_0, \dots, \hat{x}_n)_{\text{MAP}} = \underset{x_0, \dots, x_n}{\operatorname{argmax}} p(x_0, \dots, x_n | y_1, \dots, y_n) = \underset{x_0, \dots, x_n}{\operatorname{argmax}} p(x_0, \dots, x_n, y_1, \dots, y_n) \quad (8.23)$$

$$\text{Probability of the observation: } p(y_1, \dots, y_n) = \sum_{x_0, \dots, x_n} p(x_0, \dots, x_n, y_1, \dots, y_n) \quad (8.24)$$

- Note their similar forms
- Brute force is infeasible

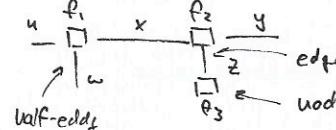
$$1) \text{ If } (\hat{x}_0, \dots, \hat{x}_n) \text{ is unique, then } \hat{x}_n = \underset{x_n}{\operatorname{argmax}} p(x_n, y_1, \dots, y_n) \triangleq \max_{x_1, \dots, x_n \text{ except } x_n} p(x_0, \dots, x_n, y_1, \dots, y_n)$$

8.2 Factor Graphs

Here: Factor graphs, Represents the factorization of a function of several vars.

Ex $f(u, w, x, y, z) = f_1(u, w, x) f_2(x, y, z) f_3(z)$

↑ global function
Factors



- there is a (unique) node for every factor
- there is a (unique) edge or half-edge for every variable
- node g is connected to edge x iff g is a function of x

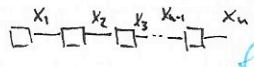
- A configuration is a particular assignment of values to all variables.
- The configuration space is the set of all configurations, i.e. the domain of the global function.
- A configuration w is called valid if $f(w) \neq 0$.
- In each configuration, the variables in the graph have a definite value. → We can consider the variables as functions of the configuration space:

$$X: (u, w, x, y, z) \mapsto x$$

Statistical models as factor graphs

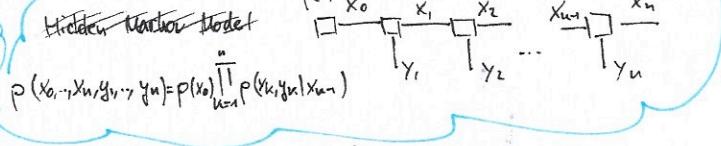
Markov Chain

$$P(x_1, \dots, x_n) = p(x_1)p(x_2|x_1)\dots p(x_n|x_{n-1})$$



Hidden Markov Model

$$P(x_0, \dots, x_n, y_1, \dots, y_n) = P(x_0) \prod_{k=1}^n P(x_k|x_{k-1}) P(y_k|x_k)$$



Equality check function

$$f_e(x, y, z) \stackrel{(8.6)}{=} \begin{cases} 1, & \text{if } x=y=z \\ 0, & \text{else} \end{cases} = \delta[x-y]\delta[y-z]$$

$$x \stackrel{f_e}{\overline{\mid}} y$$

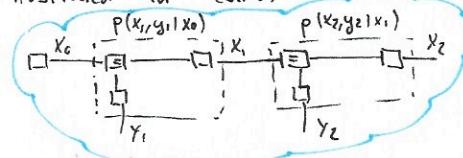
Enables this: $P(x) = f_1(x) f_2(x) f_3(x)$

$$A. \quad \square \xrightarrow{x} \square \stackrel{f_2}{\overline{\mid}} \square \xrightarrow{x' x''} \square \xrightarrow{f_3} \square$$

because $x=x'=x''$ holds for all valid configurations

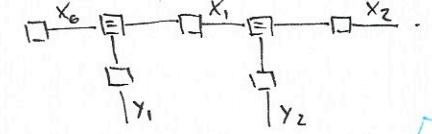
Boxes

(8.9) may be viewed as a refined version of (8.6) as illustrated in (8.10)



Hidden Markov Model

$$P(x_0, \dots, x_n, y_1, \dots, y_n) = P(x_0) \prod_{k=1}^n P(x_k|x_{k-1}) P(y_k|x_k)$$



Constraints relation such as $f_e(x, y, z) = \delta[x-y]\delta[y-z]$ is an example of a constraint. A deterministic relation such as $z = g(x, y)$ may be expressed by the factor $\delta[z - g(x, y)]$

Known & Constant Variables

Example

$$\square \xrightarrow{x_0} \square \xrightarrow{x_1} \square \xrightarrow{x_2} \dots$$

fixed observations

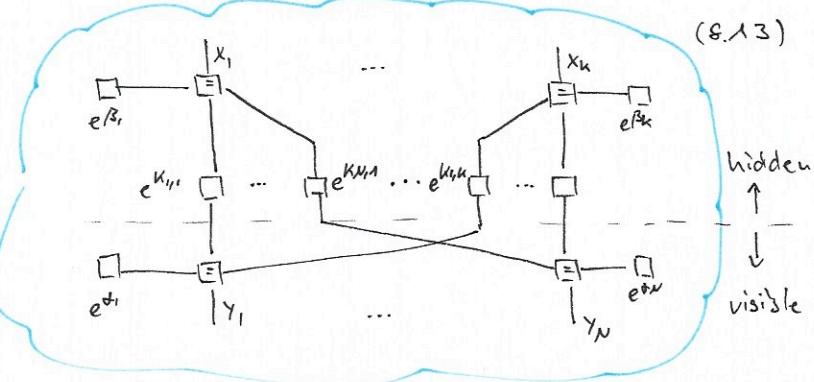
May be plugged into the factor in which case the corresponding edges can be removed.

Ex 8.2 (Restricted Boltzmann Machine)

$$P(x_1, \dots, x_n, y_1, \dots, y_n) \propto \exp \left(\sum_{i=1}^n \alpha_i(y_i) + \sum_{j=1}^n \beta_j(x_j) + \sum_{i=1}^n \sum_{j=1}^n K_{ij}(y_i, x_j) \right)$$

hidden $\in \{0, 1\}$ visible $\in \{0, 1\}$

$$\alpha_i(y_i) = \begin{cases} a_i, & y_i=1, \\ 0, & y_i=0. \end{cases} \quad \beta_j(x_j) = \begin{cases} b_j, & x_j=1, \\ 0, & x_j=0. \end{cases} \quad K_{ij}(y_i, x_j) = \begin{cases} w_{ij}, & y_i=x_j=1, \\ 0, & \text{else}. \end{cases}$$



8.3 On Factor Graphs of Statistical Models

Thm 8.1 (Independent Components)

If the factor graph of a joint prob. dist. $p(x, y)$ consists of two unconnected components, then every r.v. in one component is indep. of every r.v. in the other component.

Ex 8.4

$$\square \xrightarrow{x} \square \xrightarrow{y} \square \xrightarrow{p(y)}$$

For fixed y_1, \dots, y_n s.t. $p(y_1, \dots, y_n) \neq 0$ a FG of $p(x_1, \dots, x_n | y_1, \dots, y_n)$ is also a FG of $p(x_1, \dots, x_n | y_1, \dots, y_n)$ up to a scale factor. Proof: $p(A|B) = \frac{P(A \cap B)}{P(B)} = P(A|B)$

Ex 8.5

$$\square \xrightarrow{x_0} \square \xrightarrow{x_1} \square \xrightarrow{x_2} \dots$$

represents $p(x_0, \dots, x_n, y_1, \dots, y_n)$ and $p(x_0, \dots, x_n | y_1, \dots, y_n)$

Thm 8.3 (Markov Cut Theorem) FG of joint prob. function of r.v. $x_1, \dots, x_c, y_1, \dots, y_m$ and z_1, \dots, z_n .

If edges of y_1, \dots, y_m are removed and two graphs with only x in one and z in the other. Then for any fixed y s.t. $p(y_1, \dots, y_m) \neq 0$

$$p(x_1, \dots, x_c, z_1, \dots, z_n | y_1, \dots, y_m) = p(x_1, \dots, x_c | y_1, \dots, y_m) p(z_1, \dots, z_n | y_1, \dots, y_m)$$

In other words, conditioned on $y_1 = y_1, \dots$ any subset of x_1, \dots, x_c is independent of any subset of z_1, \dots, z_n .

Consequence of Thm 8.3: If prob. fun $p(x_{\cdot \cdot})$ of some r.v. $x_{\cdot \cdot}$ factors as $p(x_0, \dots, x_n) = g_0(x_0) \prod_{k=1}^n g_k(x_{k-1}, x_k)$ Then $x_{\cdot \cdot}$ form a Markov Chain. Moreover for a hidden Markov model with fixed observations $y_1 = y_1, \dots$ the posterior distribution $p(x_0, \dots, x_n | y_1, \dots, y_n)$ is also a Markov chain.

8.4 Closing Boxes

Usually, and by default, close a box by summing or integrating over its internal variables. The resulting closed-box function is called external function of the box.

$$\begin{array}{c} \text{Diagram of a box with internal nodes } x_2, x_3, x_4, x_5 \text{ and boundary factors } f_1, f_2. \\ \Rightarrow \text{Diagram of a single node } g \text{ with no internal structure.} \end{array} \quad g(x_1, x_2, x_4) = \sum_{x_3} \sum_{x_5} f_1(x_1, x_2, x_3) f_2(x_3, x_4, x_5)$$

↑↑ sum over all possible values of x_3/x_5

Thm 8.4 Closing an inner box within some outer box doesn't change the external function of the outer box (duh..!)

Partition Sum (PS)

"closing the entire graph" = The sum of the global function over the whole configuration space.
Closing a box in the FG does not change the partition sum.

- if all factors are $\{0, 1, 3\}$ -valued, the PS equals the number of valid configurations
- if a FG represents a prob. mass fun, its PS = 1
- if a FG represents $p(y_1, \dots, y_m, x_0, \dots, x_n)$ w.r.t. observables x_i auxiliary, then the prob. of a particular observation $y_1 = y_1, \dots$ is $p(y_1, \dots, y_n) = \sum_{x_0} \sum_{x_n} p(y_1, \dots, y_n, x_0, \dots, x_n)$ which is the PS of the FG with the y_n plugged into the factors

$$\text{Diagram showing a box with nodes } x_1, x_2, x_3, a. \text{ It is closed by factor } f(\cdot, a). \text{ The result is } f(x_1, x_2, a) = \sum_{x_3} f(x_1, x_2, x_3) f(x_3, a).$$

$$(8.23) \quad \text{Diagram showing a box with nodes } x_1, x_2, a. \text{ It is closed by factor } f(\cdot, a). \text{ The result is } f(x_1, x_2, a) = \sum_{x_3} f(x_1, x_2, x_3) f(x_3, a).$$

8.5 Matrix Multiplication and Such

$R = \mathbb{R}/\mathbb{C}/\mathbb{Z} \rightarrow$ matrix $A \in \mathbb{R}^{m \times n}$

mtx may be viewed as $\{1, \dots, m\} \times \{1, \dots, n\} \rightarrow R : (x, y) \mapsto A(x, y)$

Closing boxes opening boxes can be viewed as matrix multiplication.

matrix factorization

$$(AB)(x, z) = \sum_y A(x, y) B(y, z) \quad \leftarrow \text{mtx multiplication}$$

\rightarrow denotes row index

The identity mtx corresponds to equality check fun $f(x, y) = \begin{cases} 1 & \text{if } x=y \\ 0 & \text{else} \end{cases}$

$$(8.24) \quad \text{Diagram showing a box with nodes } x, y, z. \text{ It is closed by factor } f(x, y) f(y, z). \text{ The result is } AB.$$

(8.25) is external function of fig 8.24

For vector $v \in \mathbb{R}^n : \{1, \dots, n\} \rightarrow R$ mtx product VA

$$\text{inner product } \sum_x v(x) w(x) \quad \text{Diagram showing a box with nodes } v, w.$$

Diagonal mtx with diag vect $v \cdot v^T$

$$\text{outer product } v w^T \quad \text{Diagram showing a box with nodes } v, w.$$

Element wise multiplication (Hadamard product)

$$\text{squared vector norm } \|v\|^2 = v^T v \quad \text{Diagram showing a box with nodes } v, v.$$

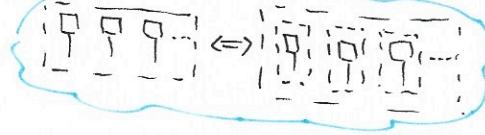
$$\text{trace} \quad \text{Diagram showing a box with nodes } A.$$

$$\text{squared Frobenius norm } \|A\|^2 = \sum_{k=1}^n \sum_{l=1}^n (A_{k,l})^2 \quad \text{Diagram showing a box with nodes } A.$$

8.6 Sum-Product Message Passing (Belief Propagation)

8.6.1 Problem Partition sum $Z_f \stackrel{(A.61)}{=} \sum_{x_1, \dots, x_n} f(x_1, \dots, x_n)$ or marginal $f_{X_k}(x_k) \stackrel{(A.62)}{=} \sum_{\substack{x_1, \dots, x_n \\ \text{except } x_k}} f(x_1, \dots, x_n)$ of many vars x_1, \dots, x_n .
 → Complexity exponential in n ! Eg if $x_i \in \{0,1\}$ (8.61) has 2^n , (8.62) 2^{n-1} terms.
 But if f can be factorized as $f(x_1, \dots, x_n) = f_1(x_1) \cdots f_n(x_n)$ (8.61) $\rightarrow Z_f = \left(\prod_{x_1} f_1(x_1) \right) / \left(\prod_{x_2} f_2(x_2) \right) \cdots$ n sums $O(n)$ linear complexity in n

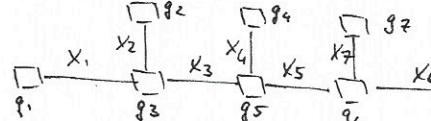
The following algorithm applies to FFG w/o cycles which then allows computing (8.61) & (8.62) with linear complexity in n



8.6.2 The Algorithm

Running example: $f_{X_3}(x_3) \stackrel{(8.61)}{=} \sum_{x_1, \dots, x_2} f(x_1, \dots, x_7)$
 $\exp. x_3$

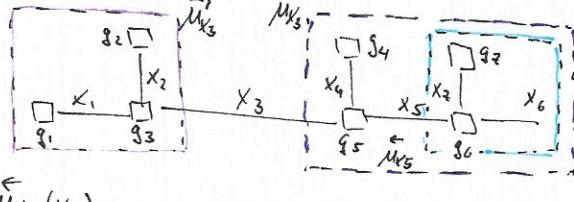
$$f(x_1, \dots, x_7) = g_1(x_1) g_2(x_2) g_3(x_1, x_2, x_3) \cdots$$



I Recursively close boxes

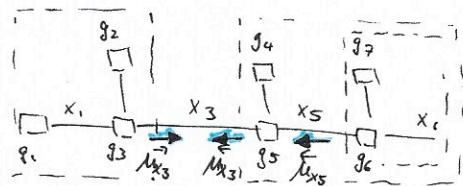
$$f_3(x_3) = \left(\sum_{x_1, x_2} g_1(x_1) g_2(x_2) g_3(x_1, x_2, x_3) \right)$$

$$\left(\sum_{x_4, x_5} g_4(x_4) g_5(x_3, x_4, x_5) \left(\sum_{x_6, x_7} g_6(x_5, x_6, x_7) g_7(x_7) \right) \right) = \overrightarrow{M}_{X_3}(x_3) \overleftarrow{M}_{X_3}(x_3)$$



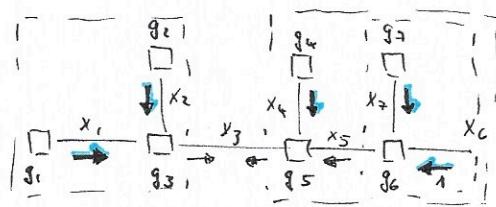
corresponding closed-box functions

II Messages View the closed-box functions \overrightarrow{M}_k as messages. These msg. are passed left to right \overrightarrow{M}_{X_3} or right to left \overleftarrow{M}_{X_3} . They are functions, e.g. $\overrightarrow{M}_{X_3}, \overleftarrow{M}_{X_3}$ are fun of x_3 .



III Messages out of leaf nodes These are simply the corresponding factors themselves: $\overrightarrow{M}_{X_1}(x_1) = g_1(x_1)$

$$\text{Half-edges: } \overleftarrow{M}_{X_6}(x_6) = 1 \quad \forall x_6$$



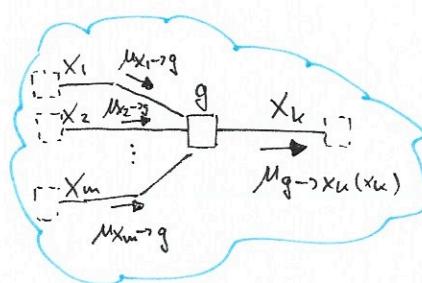
Sum-Product Message Computation Rule

The msg out of a leaf node along edge X is the function $g(x)$ itself.

The msg out of a non-leaf node along edge X_k is the function

(8.69)

$$M_{g \rightarrow X_k}(x_k) \stackrel{\#}{=} \sum_{\substack{x_1, \dots, x_m \\ \text{exc. } x_k}} g(x_1, \dots, x_m) \underbrace{M_{X_1 \rightarrow g}(x_1) \cdots M_{X_m \rightarrow g}(x_m)}_{\text{except } M_{X_k \rightarrow g}(x_k)}$$



Thm 8.5 (Sum-Product Marginals) Let $f(x_1, \dots, x_n)$ be represented by a cycle-free FFG and all messages are computed according to

$$(8.69), \text{ then } \sum_{\substack{x_1, \dots, x_n \\ \text{except } x_k}} f(x_1, \dots, x_n) = \overrightarrow{M}_{X_k}(x_k) \overleftarrow{M}_{X_k}(x_k) \quad (8.70)$$

where $\overrightarrow{M}_{X_k}/\overleftarrow{M}_{X_k}$ are the two messages along the edge X_k .

(This thm is just a consequence of thm 8.4)

8.6.3 Application to Partition Sums

From a marginal function $f_{X_k}(x_k) \stackrel{(8.62)}{=} \sum_{\substack{x_1, \dots, x_n \\ \text{except } x_k}} f(x_1, \dots, x_n)$ the partition sum is

(8.71)

$$Z_f = \sum_{X_k} f_{X_k}(x_k)$$

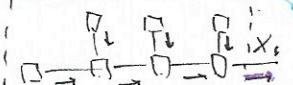
In a cycle-free FFG this

$$Z_f = \sum_{X_k} \overrightarrow{M}_{X_k}(x_k) \overleftarrow{M}_{X_k}(x_k)$$

$$Z_f = \sum_{X_1} g_1(x_1) \overleftarrow{M}_{X_1}(x_1)$$



$$Z_f = \sum_{X_6} \overrightarrow{M}_{X_6}(x_6)$$



Ex 8.6 (Counting unconsecutive binary sequences)

Sequence x_1, \dots, x_n , $x_i \in \{0,1\}$. $f: \{0,1\}^n \rightarrow \{0,1\}$
 constraint function

$$\begin{array}{c} 0101 \\ 1010 \\ 0100 \\ 0010 \end{array} \xrightarrow{\quad} \begin{array}{c} 0110 \\ 1100 \\ 1111 \\ 0011 \end{array} \times$$

$$f(x_1, \dots, x_n) = \prod_{i=1}^{n-1} K(x_i, x_{i+1})$$

$$K(x_n, x_{n+1}) = \begin{cases} 0 & x_n = x_{n+1} = 1 \\ 1 & \text{otherwise} \end{cases}$$

FG of $\sum_{n=5}^{\infty}$

$\rightarrow f(x_1, \dots, x_n)$ is only =1 if the sequence x_1, \dots, x_n has no consecutive ones.

The desired number of valid sequences is the number of valid configurations $\rightarrow =$ the partition sum Z_f

e.g. $Z_f = \sum_{x_n} \vec{\mu}_{x_n}(x_n) = \vec{\mu}_{x_1}(0) + \vec{\mu}_{x_1}(1)$. $\vec{\mu}_{x_1}(x_1) = 1$, $\vec{\mu}_{x_e}(x_e) = \sum_{x_{e-1}} K(x_{e-1}, x_e) \vec{\mu}_{x_{e-1}}(x_{e-1})$ by (8.69)

$$\vec{\mu}_{x_2}(x_2) = K(0, x_2) + K(1, x_2) = \begin{cases} 2 & x_2 = 0 \\ 1 & x_2 = 1 \end{cases}$$

$$\vec{\mu}_{x_3}(x_3) = K(0, x_3) \vec{\mu}_{x_2}(0) + K(1, x_3) \vec{\mu}_{x_2}(1) = \begin{cases} 2+1 & x_3 = 0 \\ 2+0 & x_3 = 1 \end{cases} = \begin{cases} 3 & x_3 = 0 \\ 2 & x_3 = 1 \end{cases}$$

$$\vec{\mu}_{x_4}(x_4) = K(0, x_4) \vec{\mu}_{x_3}(0) + K(1, x_4) \vec{\mu}_{x_3}(1) = \begin{cases} 5 & x_4 = 1 \\ 3 & x_4 = 0 \end{cases}$$

$$\vec{\mu}_{x_5}(x_5) = \begin{cases} 8 & x_5 = 1 \\ 5 & x_5 = 0 \end{cases}$$

$$Z_f = \sum_{x_5} \vec{\mu}_{x_5}(x_5) = 8+5 = 13 \quad \checkmark$$

8.6.4 Application to Marginal Probabilities Let $p(x_1, \dots, x_n)$ be a joint prob. density represented by a cycle-free FG.

Then $p(x_n) = \sum_{\substack{x_1, \dots, x_n \\ \text{except } x_n}} p(x_1, \dots, x_n) = \vec{\mu}_{x_n}(x_n) \vec{\mu}_{x_n}(x_n)$

can be directly computed in FG of f and then recover $\vec{\mu}_f$ from $\sum p(x_n) = 1$

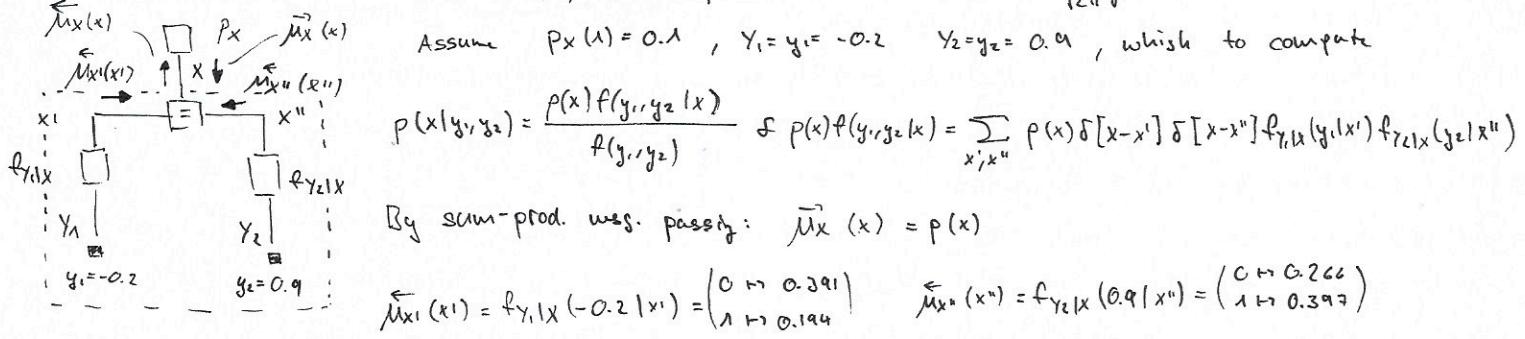
More generally a FG represents $f(x_1, \dots, x_n) = \vec{\mu}_f p(x_1, \dots, x_n)$ $\forall \neq 0$, $\vec{\mu}_f^{-1} = Z_f$ $\vec{\mu}_f^{-1} p(x_n) = \vec{\mu}_{x_n}(x_n) \vec{\mu}_{x_n}(x_n)$

(8.84)

This applies to conditional probs:

$$p(x_1 | y_1, \dots, y_m) = \sum_{\substack{x_1, \dots, x_n \\ \text{exc. } x_k}} p(x_1, \dots, x_n | y_1, \dots, y_m) \leftarrow \sum_{\substack{x_1, \dots, x_n \\ \text{exc. } x_k}} p(x_1, \dots, x_n, y_1, \dots, y_m)$$

Ex 8.8 (Noisy Indep. Measurements) $X \in \{0,1\}$, y_1, y_2 indep. meas of X $f(y_n | X) = \frac{1}{(2\pi)^{\sigma^2}} e^{-\frac{(y_n - x)^2}{2\sigma^2}}$ $n=1,2$ $\sigma=1$



Close box: $\vec{\mu}_x(x) = \sum_{x', x''} \vec{\mu}_{x_1}(x_1) \vec{\mu}_{x_2}(x_2) \delta[x - x'] \delta[x - x''] = \vec{\mu}_{x_1}(x) \vec{\mu}_{x_2}(x) = \begin{pmatrix} 0 & 0.104 \\ 1 & 0.077 \end{pmatrix}$

Finally $p(x | y_1, y_2) \leftarrow \vec{\mu}_x(x) / \vec{\mu}_x(x) = \begin{pmatrix} 0 & 0.0104 \\ 1 & 0.1 - 0.077 \end{pmatrix}$, by normalization $p(x | y_1, y_2) = \begin{pmatrix} 0 & 0.924 \\ 1 & 0.076 \end{pmatrix}$

8.7 Closing Boxes by Maximization & Max-Product Message Passing

Instead of closing boxes by sum/integrate we can do so by maximizing over all internal variables

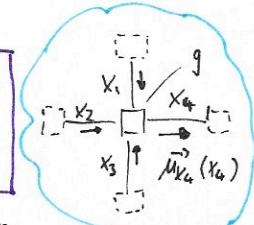
Thm 8.6 Assuming all nodes/factors are real and nonnegative, closing an inner box within an outer box by maximizing over its internal variables doesn't change the external function of the outer box.

It follows that message passing still holds when replacing summation with maximization

Max-Product Message Computation Rule

$$Mg \rightarrow x_n(x_n) = \vec{\mu}_{x_n}(x_n) = \begin{cases} g(x_n) & \text{if leaf} \\ \max_{\text{nonleaf}} & \text{else} \end{cases}$$

$$Mg \rightarrow x_n(x_n) \triangleq \max_{\substack{x_1, \dots, x_m \\ \text{except } x_n}} g(x_1, \dots, x_m) \underbrace{\vec{\mu}_{x_1 \rightarrow g}(x_1) \cdots \vec{\mu}_{x_m \rightarrow g}(x_m)}_{\text{except } Mg \rightarrow g(x_n)}$$



Thm 8.7 (Max-Product Marginals)

Let $f(x_1, \dots, x_n)$ represented by cycle-free FG with real and nonnegative factors/nodes then

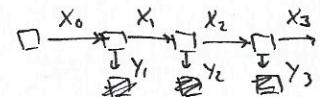
$$\max_{\substack{x_1, \dots, x_n \\ \text{except } x_n}} f(x_1, \dots, x_n) = \vec{\mu}_{x_n}(x_n) \vec{\mu}_{x_n}(x_n)$$

Importance:
 $(x_1, \dots, x_n) \triangleq \arg\max_{x_1, \dots, x_n} f(x_1, \dots, x_n)$
 $\hat{x}_n = \arg\max_{x_n} \vec{\mu}_{x_n}(x_n) \vec{\mu}_{x_n}(x_n)$

8.8 Message Passing in Markov Chains and Hidden Markov Models (HMM)

8.8.1 Msg Passing in HMM

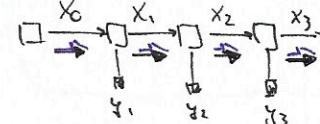
Recall characterization of HMM: $p(x_0, \dots, x_n, y_1, \dots, y_n) = p(x_0) \prod_{k=1}^n p(x_k, y_k | x_{k-1})$



Current state estimation

$$p(x_n | y_1, \dots, y_n) \triangleq \sum_{x_0, \dots, x_{n-1}} p(x_0, \dots, x_n, y_1, \dots, y_n) = \overrightarrow{M}_{X_n}(x_n)$$

can be calculated by forward sum-prod. msg pass



Prediction of next output:

$$p(y_n | y_1, \dots, y_{n-1}) \triangleq \sum_{x_0, \dots, x_{n-1}} p(x_0, \dots, x_n, y_1, \dots, y_n) = \overrightarrow{M}_{Y_n}(y_n)$$

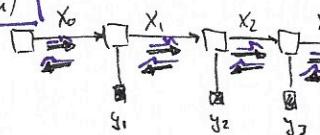
forward sum-prod. msg pass



Estimation of state at any time:

$$p(x_k | y_1, \dots, y_n) \triangleq \sum_{\substack{x_0, \dots, x_{k-1} \\ \text{except } x_k}} p(x_0, \dots, x_n, y_1, \dots, y_n) = \overrightarrow{M}_{X_k}(x_k) / \overleftarrow{M}_{X_k}(x_k)$$

forward-backward sum-prod. msg. passif



MAP estimate of state trajectory:

$$(x_0, \dots, x_n)_{\text{MAP}} \stackrel{?}{=} \underset{x_0, \dots, x_n}{\text{argmax}} p(x_0, \dots, x_n | y_1, \dots, y_n) \quad \text{if } \text{argmax} \text{ is unique}$$

$$\rightarrow \hat{x}_k = \underset{x_k}{\text{argmax}} \overrightarrow{M}_{X_k}(x_k) / \overleftarrow{M}_{X_k}(x_k)$$

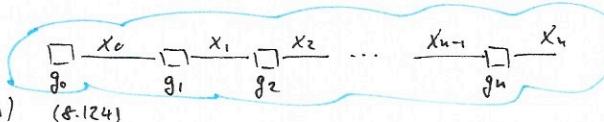
Scale factors in \overrightarrow{M}_k cannot be ignored in this case!

Probability of observation

$$\begin{aligned} p(y_1, \dots, y_n) &= \sum_{x_0, \dots, x_n} p(x_0, \dots, x_n, y_1, \dots, y_n) \\ &= \sum_{x_n} \sum_{x_0, \dots, x_{n-1}} p(x_0, \dots, x_n, y_1, \dots, y_n) = \sum_{x_n} \overrightarrow{M}_{Y_n}(x_n) \end{aligned}$$

8.8.2 MC Sampling and Reparameterization

If x_0, \dots, x_n is a MC $p(x_0, \dots, x_n) = p(x_0)p(x_1|x_0)p(x_2|x_1) \dots p(x_n|x_{n-1})$



Creating a sample (x_0, \dots, x_n) from (8.124) (i.e. simulating the MC) is straightforward:

1. draw x_0 from $p(x_0)$
2. draw x_1 from $p(x_1 | x_0)$
3. etc... (8.125)

Sampling is less obvious, but if $p(x_0, \dots, x_n)$ is given by $\overrightarrow{g_0(x_0)} \prod_{k=1}^n g_k(x_{k-1}, x_k)$ with general factors g_k

Ex 8.9 (Nondegenerate Sim seq.): Wish to generate rand seq $(x_0, \dots, x_n) \in \{0, 1\}^{n+1}$ w/o consecutive ones, uniformly dist. over all such seq. This is a MC with $g_0(x_0) = 1$ and $g_k(x_{k-1}, x_k) = \begin{cases} 0 & \text{if } x_{k-1} = x_k = 1 \\ 1 & \text{else} \end{cases}$

Ex 8.11 (Sampling from MC with fixed final state): x_0, x_1, \dots, x_n a MC with fixed $x_n = x_0$. Wish to

draw (x_0, \dots, x_n) according to $p(x_0, \dots, x_{n-1} | x_n)$.

Since (8.125) is MC can be parametrized as (8.124) st. $p(x_0) = g_0(x_0) / \overrightarrow{M}_{X_0}(x_0)$

$$\text{and } \overrightarrow{p(x_n | x_{n-1})} = \frac{p(x_n, x_{n-1})}{p(x_{n-1})} = \frac{\overrightarrow{M}_{X_{n-1}}(x_{n-1}) g_n(x_{n-1}, x_n) / \overleftarrow{M}_{X_n}(x_n)}{\overrightarrow{M}_{X_{n-1}}(x_{n-1}) \overleftarrow{M}_{X_{n-1}}(x_{n-1})} = \frac{\overrightarrow{g_n(x_{n-1}, x_n)} / \overleftarrow{M}_{X_n}(x_n)}{\overrightarrow{M}_{X_{n-1}}(x_{n-1})} \quad (\text{8.130})$$

Sampling can thus be done by backwards sum-prod. msg pass followed by forward sample.

8.9 Additional Remarks and Details

Message Scaling $\tilde{\mu}(\cdot)$ typically tend to quickly zero/bifurcate. In practice they are scaled/normalized. May have to keep track of scale factor, in log form.

Message Representation $\tilde{\mu}(\cdot)$ are functions. If domain is finite, store tables. If domain is binary, store $\frac{\tilde{\mu}(0) - \tilde{\mu}(1)}{\tilde{\mu}(0) + \tilde{\mu}(1)}$ or $\log \frac{\tilde{\mu}(0)}{\tilde{\mu}(1)}$

Msg. passing through eq. constraints

$$\begin{array}{c} X \\ \xrightarrow{\quad} \boxed{=} \xrightarrow{\quad} Z \\ Y \downarrow \end{array} \quad \tilde{\mu}_Z(z) = \sum_x \sum_y f_{x,y,z}(x,y,z) \tilde{\mu}_X(x) \tilde{\mu}_Y(y) = \tilde{\mu}_X(z) \tilde{\mu}_Y(z)$$

or $\tilde{\mu}_Z(z) = \max \dots = \tilde{\mu}_X(z) \tilde{\mu}_Y(z)$
for max-prod.

Backward Msg in Chain-Rule Models FG of $p(x,y) = p(x)p(y|x)$

$$\begin{array}{c} X \\ \xrightarrow{\quad} \boxed{=} \xrightarrow{\quad} Y \\ P_X \quad P_{Y|X} \end{array} \quad \tilde{\mu}_X(x) = 1$$

Y unknown

$$\begin{array}{c} X \\ \xleftarrow{\quad} \boxed{=} \xleftarrow{\quad} Y \\ P_X \quad P_{Y|X} \end{array} \quad \tilde{\mu}_X(x) = P_{Y|X}(y|x)$$

Y known

likelihood function

Message passing in FG with cycles

- Result can be arbitrarily bad
- Is used with excellent rates, e.g. in decoding error correcting codes
- Initialize all $\tilde{\mu} = \tilde{\mu}^* = 1$
- Iterate through nodes with some schedule
- Stop after time / condition

Chapter 10: Gaussian Message Passing

10.1 Introduction Notes: • Generalization of sum-prod. msg passing to continuous variables/func. is by replacing sum by integral $\delta[\cdot] \rightarrow \delta(\cdot)$

Ex10.1 (cf. Ex8.8)

Given: $X \in \mathbb{R}$ r.v. $f(x)$, y_1, y_2 obsrvations with $f(y_1, y_2 | x) = f(y_1 | x) f(y_2 | x)$ $y_1 = y_1$, $y_2 = y_2$

Problem: Find posterior $f(x | y_1, y_2)$

$$\text{Sol: } f(x | y_1, y_2) \propto \tilde{\mu}_X(x) \tilde{\mu}_{y_1}(x) \tilde{\mu}_{y_2}(x) \quad \tilde{\mu}_X(x) = f(x) \quad \tilde{\mu}_{y_1}(x) = f_{y_1}(y_1 | x) \quad \tilde{\mu}_{y_2}(x) = f_{y_2}(y_2 | x)$$

$$\tilde{\mu}_X(x) = \int_{x_1} \int_{x_2} \tilde{\mu}_X(x_1) \tilde{\mu}_X(x_2) \delta(x-x_1) \delta(x-x_2) dx_1 dx_2 = f_{y_1}(y_1 | x) f_{y_2}(y_2 | x)$$

Focus lies now on scalar Gaussian messages described by mean m and variance σ^2

$$\begin{array}{c} X \xrightarrow{\quad} Z \\ Y \downarrow \\ \delta(x-y) \delta(x-z) \end{array} \quad \begin{aligned} \tilde{\mu}_Z &= \frac{\tilde{\mu}_X/\tilde{\sigma}_X^2 + \tilde{\mu}_Y/\tilde{\sigma}_Y^2}{1/\tilde{\sigma}_X^2 + 1/\tilde{\sigma}_Y^2} \quad (10.11) \\ 1/\tilde{\sigma}_Z^2 &= 1/\tilde{\sigma}_X^2 + 1/\tilde{\sigma}_Y^2 \quad (10.12) \\ \tilde{\mu}_X &= \frac{\tilde{\mu}_Y/\tilde{\sigma}_Y^2 + \tilde{\mu}_Z/\tilde{\sigma}_Z^2}{1/\tilde{\sigma}_Y^2 + 1/\tilde{\sigma}_Z^2} \quad 1/\tilde{\sigma}_X^2 = 1/\tilde{\sigma}_Y^2 + 1/\tilde{\sigma}_Z^2 \end{aligned}$$

$$m(x) \propto \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right) \propto \exp\left(-x^2 \frac{1}{2\sigma^2} + x \frac{m}{\sigma^2}\right)$$

$$\begin{aligned} \tilde{\mu}_Z &= \tilde{\mu}_X + \tilde{\mu}_Y \\ \tilde{\sigma}_Z^2 &= \tilde{\sigma}_X^2 + \tilde{\sigma}_Y^2 \\ \tilde{\mu}_X &= -\tilde{\mu}_Y + \tilde{\mu}_Z \\ \tilde{\sigma}_X^2 &= \tilde{\sigma}_Y^2 + \tilde{\sigma}_Z^2 \end{aligned}$$

$$\begin{aligned} X &\xrightarrow{\quad} Y \\ \delta(y-ax) & \end{aligned}$$

$$\begin{aligned} \tilde{\mu}_Y &= a \tilde{\mu}_X \\ \tilde{\sigma}_Y^2 &= a^2 \tilde{\sigma}_X^2 \\ \tilde{\mu}_X &= \tilde{\mu}_Y/a \\ \tilde{\sigma}_X^2 &= \tilde{\sigma}_Y^2/a^2 \end{aligned}$$

Incoming Gaussian messages to these nodes are also Gaussian at their outputs.

Note: The means never affect the variances (but not vice-versa)

Ex10.2 (Ex10.1 contd) Let $f(x) = \mathcal{N}(m_x, \sigma_x^2)$ $f(y_1|x) = \mathcal{N}(x, \sigma^2)$. We then have $\tilde{\mu}_X = m_x$ $\tilde{\sigma}_X^2 = \sigma^2$ $\tilde{\mu}_{y_1} = y_1$ $\tilde{\sigma}_{y_1}^2 = \sigma^2$

$$\text{Then by (10.11)/(10.12)} \quad \tilde{\mu}_X = (y_1 + y_2)/2 \quad \tilde{\sigma}_X^2 = \sigma^2/2 \quad \tilde{\mu}_{y_1y_2} = \frac{\tilde{\mu}_X/\tilde{\sigma}_X^2 + \tilde{\mu}_{y_2}/\tilde{\sigma}_{y_2}^2}{1/\tilde{\sigma}_X^2 + 1/\tilde{\sigma}_{y_2}^2} = \frac{m_x/\sigma^2 + (y_2 + \sigma^2)/\sigma^2}{1/\sigma^2 + 2/\sigma^2} \quad \tilde{\sigma}_{y_1y_2}^2 = 1/\sigma^2 + 2/\sigma^2$$

Remarks from Appendix C Facts about Gaussian vectors:

1. For Gaussian msg., sum-prod. coincides with max-prod. and passivity
2. For jointly Gaussian r.v.: $MAP = MMSE = LMSE$ estimates
3. Even if X, Y not jointly Gaussian, the LMSE of X from $Y=y$ may be obtained by pretending X, Y jointly Gaussian and having their corresponding MAP

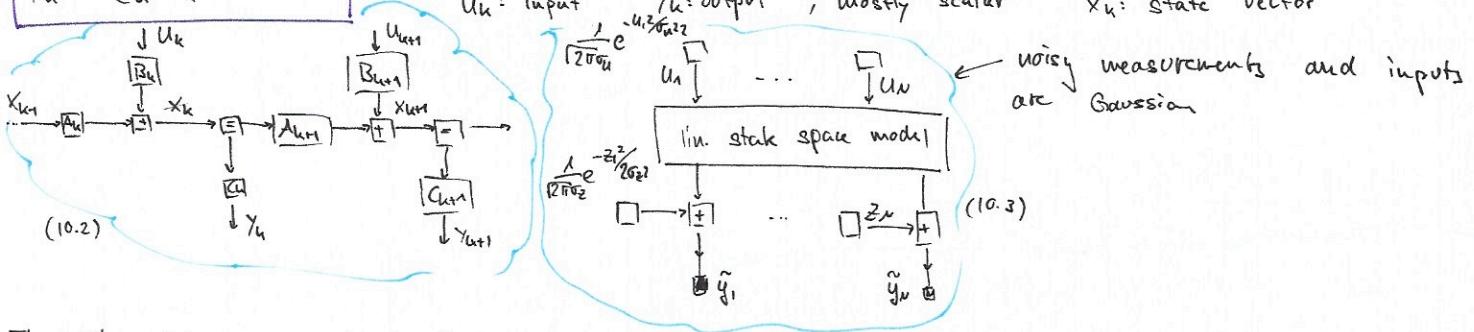
Merge: $m = \frac{\tilde{m}/\tilde{\sigma}^2 + \tilde{m}/\tilde{\sigma}^2}{1/\tilde{\sigma}^2 + 1/\tilde{\sigma}^2}$ $\sigma^2 = \frac{1}{1/\tilde{\sigma}^2 + 1/\tilde{\sigma}^2}$

10.2 Linear Gaussian State Space Models

$$\begin{aligned} X_k &= A_k X_{k-1} + B_k U_k \\ Y_k &= C_k X_k \end{aligned}$$

Linear state space model:

U_k, X_k, Y_k real col. vectors A_k, B_k, C_k real matrices
 U_k : input Y_k : output, mostly scalar X_k : state vector



The initial state X_0 and input signal U_1, \dots, U_N fully determine both the state trajectory X_1, \dots, X_N and the noise free output Y_1, \dots, Y_N . In consequence, the MAP/MME/MMSE of any of these quantities is fully determined by the minimizer of

$$\underset{x_0, u_1, \dots, u_N}{\operatorname{argmax}} \prod_{k=1}^N e^{-\frac{U_k^2}{2\sigma_u^2}} e^{-\frac{(y_k - \tilde{y}_k)^2}{2\sigma_z^2}} = \underset{x_0, u_1, \dots, u_N}{\operatorname{argmin}} \left(\frac{1}{\sigma_u^2} \sum_{k=1}^N U_k^2 + \frac{1}{\sigma_z^2} \sum_{k=1}^N (y_k - \tilde{y}_k)^2 \right) \quad (10.3)$$

10.3 Vector-Gaussian Message Passing and Kalman Filtering

Gaussian Message Passing through graph (10.3) with (10.2) plugged in is known as Kalman Filter!

Recall: Vector \vec{x} of Gaussian r.v.
 $\xrightarrow{\text{real}}$

$$f_x(\vec{x}) \propto e^{-\frac{1}{2}(\vec{x} - \vec{m})^\top \vec{V}^{-1} (\vec{x} - \vec{m})}$$

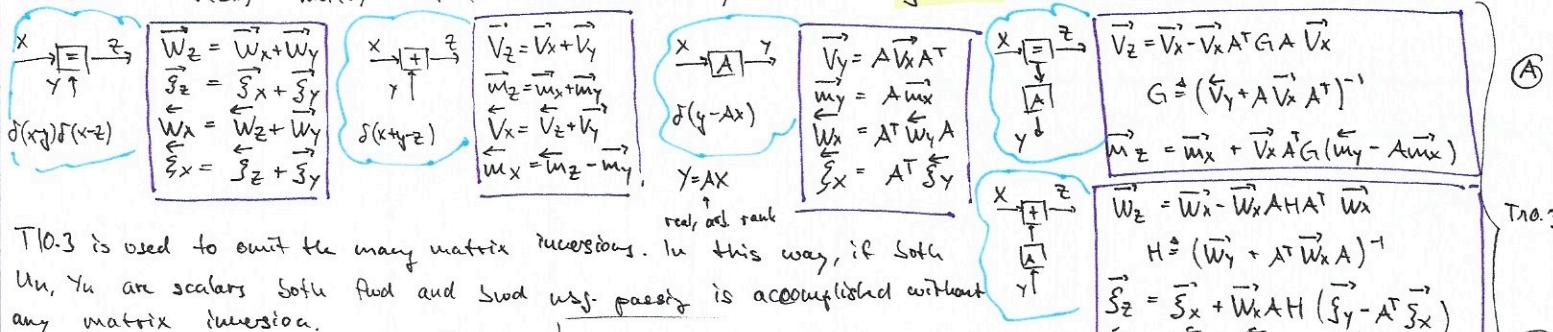
\vec{m} : mean vector

\vec{V} : covariance matrix

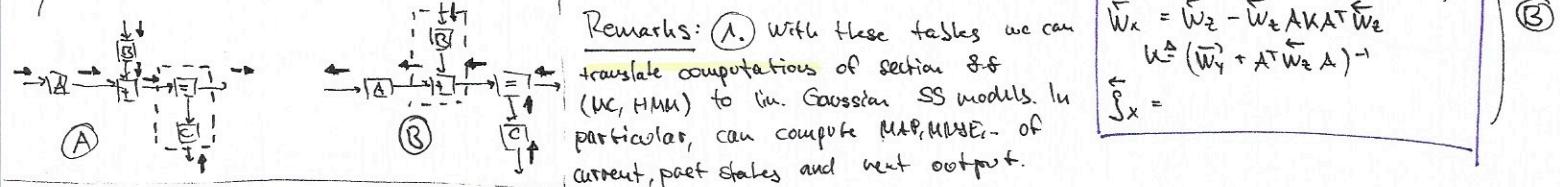
$\vec{W} = \vec{V}^{-1}$: precision matrix

$\vec{s} \triangleq \vec{W}\vec{m}$: scale vector

$$\text{exact: } f_x(\vec{x}) = \frac{1}{\sqrt{(2\pi)^n}} \cdot \frac{1}{\det(V)} \exp\left(-\frac{1}{2}(\vec{x} - \vec{m})^\top \vec{V}^{-1} (\vec{x} - \vec{m})\right)$$

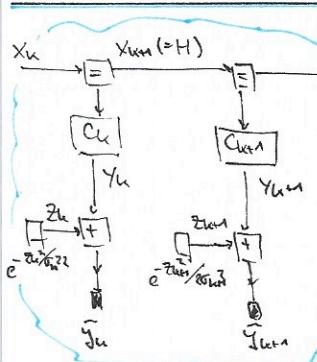


T10.3 is used to omit the many matrix inversions. In this way, if both U_k, Y_k are scalars both fwd and bwd msg passing is accomplished without any matrix inversion.



- ② The final (MAP,...) estimate of some quantity X is the mean vector m_X of the scaled Gaussian $\mu_X(x) \triangleq \vec{\mu}_X(x) \vec{\mu}_X(x)$
- ③ A known variable s.a. $\tilde{y}_k = y_k$ leads to degenerate Gaussian $\mu_{\tilde{y}_k}$, $\vec{\mu}_{\tilde{y}_k} = \tilde{y}_k$, $\vec{\nu}_{\tilde{y}_k} = 0$. If not observed $\vec{\mu}_{\tilde{y}_k} = 0$ $\vec{\nu}_{\tilde{y}_k} = 0$
- ④ Comp. rules for \vec{V}, \vec{W} do not involve any mean vectors \rightarrow can be precomputed before \tilde{y}_k is observed.
 Applies in particular in stationary where A_k, B_k, C_k do not depend on k
- ⑤ Numerical problems will arise. Writing $A_k = I$ in Jordan form often helps

10.5 Recursive Least Squares (RLS) and Adaptive Filters



(10.6) may be viewed as (10.2)/(10.3) with $A_k = I$, $B_k = 0$. X_k : unknown real-valued col. vector, C_1, C_2, \dots real valued row vectors. The unknown $X_k = X_{k+1} = \dots = H$

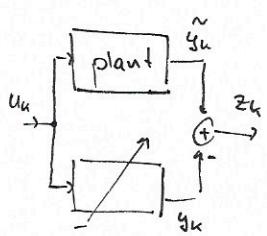
(10.6) expresses $\tilde{y}_k = C_k H + z_k$, z_1, \dots real, zero-mean, indep. Gaussian.

10.5.1 RLS Recall LS-problem with $\|Ah - \tilde{y}\|^2$ find h . Let C_h be h -th component of \tilde{y} . Then $\|Ah - \tilde{y}\|^2 = \sum_x (C_h h - \tilde{y}_h)^2$ thus

$$\underset{h}{\operatorname{argmin}} \|Ah - \tilde{y}\|^2 = \underset{h}{\operatorname{argmax}} e^{-\frac{1}{2} \|Ah - \tilde{y}\|^2} = \underset{h}{\operatorname{argmax}} \prod_k e^{-(C_h h - \tilde{y}_h)^2}$$

$\mathbb{A} = \prod_k e^{-2\pi z_k^2}$ is represented by fig 10.6 $\rightarrow \hat{h} = \hat{x}$ may be found by max- (=sum-) prod. Gaussian msg. passing

10.5.2 Adaptive FIR Filters



Given: real input u_1, u_2, \dots, u_N
 Which: final filter taps h_0, \dots, h_M s.t.
 desired real output $\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_N$
 $y_N = \sum_{k=0}^M h_k u_{N-k}$ is a good approx. of \tilde{y}_N

Specifically $\sum_{k=1}^N (\tilde{y}_k - y_k)^2$ as small as possible.

With $h = (h_0, h_1, \dots, h_M)^T$ $c_h = (u_1, u_2, \dots, u_M)$ $y_h = \sum_k h_k u_{N-k} = c_h h$

and the desired sol.: $\hat{h} = \arg \min_h \sum_k (\tilde{y}_k - c_h h)^2 = \arg \max_h \prod_k e^{-(\tilde{y}_k - c_h h)^2}$ which is again the maximization of FG 10.6

10.5.4 RLS Algorithms as Message Passing

- Actual sum-prod. = max-prod. msg. passing through 10.6 is equivalent to classic RLS algs.
- Forward-only (left to right) is enough: if the last observation is \tilde{y}_N , the desired sol. \hat{h} is the mean vector $\vec{m}_{X_{N+1}}$ of the message $\vec{M}_{X_{N+1}}$.
- Two methods: I. Propagate $\vec{W}_{X_h}, \vec{s}_{X_h}$, extract mean vector by solving $\vec{s}_{X_h} = \vec{W}_{X_h} \vec{m}_{X_h}$
 II. Propagate $\vec{V}_{X_h}, \vec{m}_{X_h}$, \vec{m}_{X_h} is already available at all times
- In every step multiply V_{X_h} in each step by $r > 1, r \approx 1$ scalar to slowly forget the past and improve numerical stability.

Chapter 12: Monte Carlo Methods (MC)

A probability density can be approx. represented by a list of samples drawn from the density. $x^{(1)}, \dots, x^{(L)}$ are (suff. indep.) samples from f_x and g a function defined on the range of X

$$E[g(x)] = \int_x g(x) f_x(x) dx \approx \frac{1}{L} \sum_{l=1}^L g(x^{(l)}) \quad (12.2)$$

with identity for $L \rightarrow \infty$ w.p.1

If g is indicator function of some set S $g(x) = \begin{cases} 1 & x \in S \\ 0 & x \notin S \end{cases}$ then

$$P(x \in S) = \int_x g(x) f_x(x) dx \approx \frac{1}{L} \sum_{l=1}^L g(x^{(l)}) = \frac{1}{L} \sum_{x \in S} 1$$

Generally it is sometimes attractive to work with samples not from f_x but from some aux. dist. g .

$$\hat{E}[g(x)] \triangleq \frac{1}{L} \sum_{l=1}^L \frac{f_x(x^{(l)})}{g(x^{(l)})} \cdot g(x^{(l)}) \quad (12.7)$$

↑ with respect to g

this is called importance sampling
 bcs the var of (12.7) may be larger/smaller than the one of (12.2)

Example: Samples $(x^{(1)}, y^{(1)}), \dots$ from $f_{x,y}$ can be turned into a list of samples from f_y by dropping the $x^{(1)}$ component of each sample.

A fundamental limitation of MC is the prob. of very rare events cannot be assessed from a list of samples of practical size. the convergence of (12.2) is not guaranteed.

12.1 Particle Filter Problem: Implementation of sum-prod. msg. passing if not Gaussian (=continuous alphabet) is hard!
 → MC can help: lists of samples may be used to represent joint dist., messages and marginals.

= seq. MC algorithm
 Illustrate by application to state estimation: Fwd. sum-prod. msg. passing 1. init: list of indep. samples $\{x_0^{(e)}\}$ from f_{X_0} represent \vec{m}_{X_0}

with samples (particles): $f(x_0, \dots, x_n, y_1, \dots, y_n) = \prod_{k=1}^n f(x_k | x_{k-1}) f(y_k | x_{k-1}, \dots, x_0)$

each sample from the previous list $\{x_n^{(e)}\}$ (with repetitions) with probability proportional to

$f_{Y_{n+1}|X_n}(y_{n+1} | x_n^{(e)})$. The new sample represent $\vec{m}_{X_{n+1}}(x) = \vec{m}_{X_n}(x) f_{Y_{n+1}|X_n}(y_{n+1} | x)$

4. The msg. $\vec{m}_{X_{n+1}}$ is represented by samples $\{x_{n+1}^{(e)}\}$, which are sampled indep. from $f_{X_{n+1}|X_n}(\cdot | x_n^{(e)})$

12.2 Gibbs Sampling General method for creating samples from some joint pros. $p(x) = p(x_1, \dots, x_n)$

* Can be used for FG with cycles → Samples are not indep., dependence b/w $x^{(e)}, x^{(e+m)}$ may decay slowly in consider case $n=2$. Samples form a Markov Chain.

1. Begin with sample $x_2^{(0)}$ from $p^{(x_2)}$
 2. $x_1^{(e)}$ is sampled according to $P_{x_1|x_2}(x_1|x_2^{(e-1)}) \leftarrow P_{x_1,x_2}(x_1, x_2^{(e-1)})$
 3. $x_2^{(e)}$ → $P_{x_2|x_1}(x_2|x_1^{(e)}) \leftarrow P_{x_1,x_2}(x_1^{(e)}, x_2)$
- } repeat

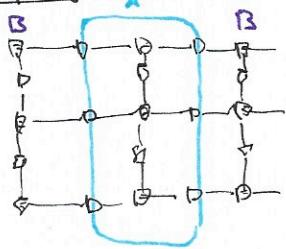
Step 1. is difficult. Practice: choose $x_2^{(0)}$ from oftr dist. or deterministic with $P_{x_2}(x_2^{(e)}) > 0$

→ need some time to converge to P_{x_1,x_2}

For $n > 2$: initial $x^{(1)}, x^{(2)}$ with for $k=1$ to n }

$$\text{sample } x_k^{(e)} \text{ from } q_k^{(e)}(x_k) = P_{x_k|x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_n}(x_k|x_1^{(e)}, \dots, x_{k-1}^{(e)}, x_{k+1}^{(e)}, \dots, x_n^{(e)})$$

Improving Gibbs Sampling



Let $p(x_1, \dots, x_n)$ be represented by FG.

→ Partition FG into sets without cycles

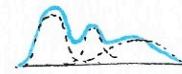
- Apply Gibbs on block level by alternating between sampling from $p(x_k|x_{\bar{k}})$ with x_B fixed and $p(x_B|x_k)$ with x_k fixed. Both steps can be carried out by backward filtering, forward sampling.

→ Example application on Restricted Boltzmann Machine

Chapter 13: Parameter Esti. w/ Expectation Maximization Goal: Determine statistical model parameters e.g. $\hat{\theta} = \arg \max \rho(y^{(1)}, \dots, y^{(L)} | \theta)$ with data $y^{(1)}$

13.1 Gaussian-Mixture Estimation (GM)

A scalar GM is a pros. density $f(y) = \sum_{k=1}^K p_k f_k(y)$ (13.2)



$$K \geq 1 \quad \sum_{k=1}^K p_k = 1 \quad f_k(y) \sim \mathcal{N}(\mu_k, \sigma_k^2), \quad \text{same as (13.2) (13.4)}$$

$$f(y) = \sum_{k=1}^K p_k f_k(y) \quad X: r.v. \in \{1, \dots, K\}$$

13.1.1 Problem Given indep. samples $y^{(1)} \in \mathbb{R}^L$ drawn from (13.2), find p_k, μ_k, σ_k^2 .

Special case $K=1$: $\hat{\mu} = \frac{1}{L} \sum_{l=1}^L y^{(l)}$, $\hat{\sigma}^2 = \frac{1}{L} \sum_{l=1}^L (y^{(l)} - \hat{\mu})^2 = \frac{1}{L} \sum_{l=1}^L (y^{(l)} - \hat{\mu})^2$ (13.5) or $\hat{\sigma}^2 = \frac{1}{L-1} \sum_{l=1}^L (y^{(l)} - \hat{\mu})^2 = \frac{1}{L-1} \sum_{l=1}^L (y^{(l)} - \hat{\mu})^2 - \frac{1}{L-1} \hat{\mu}^2$ (13.6) $\hat{\mu}, \hat{\sigma}^2$ are ML estimates

13.1.2 Algorithm For $K \geq 1$ → a closed form sol.

1. Initial guess $\hat{\theta} = (\hat{p}_1, \dots, \hat{p}_K, \hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\sigma}_1^2, \dots, \hat{\sigma}_K^2)$ (13.10)

2. For $e=1, \dots, L$, $k=1, \dots, K$ compute posterior according to current model $P(X^{(e)}=k | y^{(e)}, \hat{\theta}) = \frac{\hat{p}_k f_k(y^{(e)})}{\sum_{k=1}^K \hat{p}_k f_k(y^{(e)})}$ (13.11) $X^{(e)}$: indexes components as in (13.4)
 $\hat{f}_k(y)$: $\sim \mathcal{N}(\hat{\mu}_k, \hat{\sigma}_k^2)$

3. Compute new estimates

$$\hat{p}_k = \frac{1}{L} \sum_{l=1}^L P(X^{(l)}=k | y^{(l)}, \hat{\theta}) \quad \hat{\mu}_k = \frac{1}{\hat{p}_k} \sum_{l=1}^L y^{(l)} P(X^{(l)}=k | y^{(l)}, \hat{\theta}) \quad \hat{\sigma}_k^2 = \frac{1}{\hat{p}_k} \sum_{l=1}^L (y^{(l)} - \hat{\mu}_k)^2 P(X^{(l)}=k | y^{(l)}, \hat{\theta}) \quad (13.12)$$

$S_k \stackrel{\Delta}{=} \sum_{l=1}^L P(X^{(l)}=k | y^{(l)}, \hat{\theta})$ (13.12) is a weighted avg. of samples $y^{(l)}$ weighted by $P(X^{(l)}=k | y^{(l)}, \hat{\theta})$

4. Repeat 2,3. until convergence

13.1.3 Multivariate Case & Clustering

Generalization to Multivar.: replace (13.13) by an esti. of cov. matx. with some regularization (B.25) $P = t(L)$ e.g.

Then it can be viewed as improved k-means clustering: the "hard" assignments to classes is replaced by pds (13.10) and the isotropic Euclidean distance by a general covariance matrix.

The algo in 13.1.2 can be extended with a final clustering step that assigns classes samples to classes ("distributions"?). Can handle many cases where k-means fails.

$$\hat{V}_Y = V_S + \frac{\rho}{L} \cdot I$$

END