# ■ GENERATING OF SYNTHETIC DATA − Cont'd.
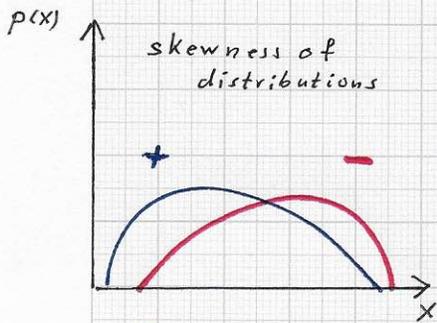
- **Review − Moments & other characteristics:**

A variety of statistical properties of feature point/vector distributions can be used to establish meaningful definitions for the difference of two distributions or error of an approximation of a distribution. Several properties are based on MOMENTS; important ones are the following (defined for discrete, finite data sets $\{x_i\}_{i=1}^{N}$):

skewness of distributions

blue p(x): + skewness
red p(x): − skewness

kurtosis of distributions

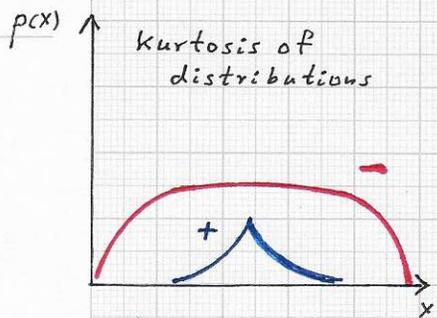blue p(x): + kurtosis
red p(x): − kurtosis

- **mean:** $\quad \bar{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$

- **variance:** $\quad \sigma^2 = \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \bar{x})^2$

- **standard deviation:** $\sigma = \sqrt{\sigma^2}$

- **skewness:** $\quad \frac{1}{N} \sum_{i=1}^{N} \left( (x_i - \bar{x})/\sigma \right)^3$

- **kurtosis:** $\quad -3 + \frac{1}{N} \sum_{i=1}^{N} \left( (x_i - \bar{x})/\sigma \right)^4$
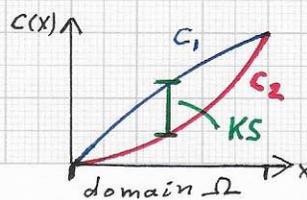
- **chi-square:** $\quad \chi^2 = \sum_{j=1}^{B} (N_j - n_j)^2 / n_j \ , \ n_j \neq 0$

( binned data, bins $1 \ldots B$; $n_j$ = true, expected no., $N_j$ = observed no. )

chi-square difference

$x_{min}$   bin j   $x_{max}$

- true, expected
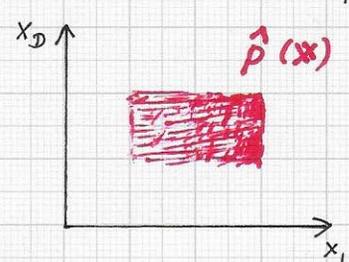- measured, observed

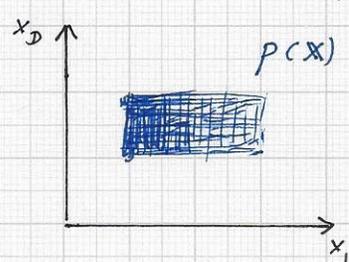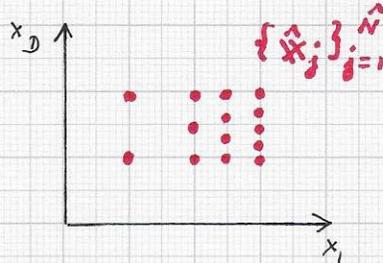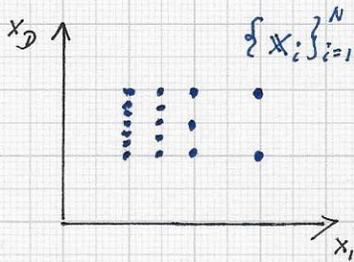- **Kolmogorov-Smirnov (KS):** $\quad KS = \max_{\Omega} |c_1(x) - c_2(x)|$

domain $\Omega$

( analytically defined data; consider absolute difference of cumulative distributions $c_1(x)$ and $c_2(x)$ )

# ■ GENERATION OF SYNTHETIC DATA - Cont'd.

• Review -Difference measures for distributions:



$\{x_i\}_{i=1}^{N}$

$\{\hat{x}_j\}_{j=1}^{\hat{N}}$

$p(x)$

$\hat{p}(x)$

given feature point samples $\{x_i\}$ and $\{\hat{x}_j\}$; constructed distribution functions $p(x)$ and $\hat{p}(x)$, used to determine similarity of $\{x_i\}$ and $\{\hat{x}_j\}$.

Functions $p(x)$ and $\hat{p}(x)$ should be normalized over their common domain $\Omega$ with volume $|\Omega|$.

One must often define and compute meaningful distances / differences between two distributions of feature point/vector data over some (finite) domain $\Omega$ in $D$-dimensional feature space. For example, one might have to compare two discrete, finite feature value distribution sets $\{x_i\}_{i=1}^{N}$ and $\{\hat{x}_j\}_{j=1}^{\hat{N}}$. The two sets each "imply" an underlying analytical function $p, \hat{p}$ that represents the respective sample $\{x_i\}, \{\hat{x}_j\}$. (Piecewise polynomial best approximations could be thought of.) Regardless of the specific method chosen to construct a function $p(x)$ from $\{x_i\}$, and $\hat{p}(x)$ from $\{\hat{x}_j\}$, one can consider a variety of measures for the comparison of $\{x_i\}$ and $\{\hat{x}_j\}$. When evaluating $p(x)$ in a uniformly integrated sense, point sets of "class $\{x_i\}$" result; similarly $\hat{p}(x)$ produces point sets of "class $\{\hat{x}_j\}$" when evaluated uniformly in integrated fashion. Thus, one can define and compute difference measures for $p$ and $\hat{p}$ to define difference of $\{x_i\}, \{\hat{x}_j\}$:

• root-mean-square (RMS): $\left( \frac{1}{|\Omega|} \int_{\Omega} \left( p(x) - \hat{p}(x) \right)^2 dx \right)^{1/2}$

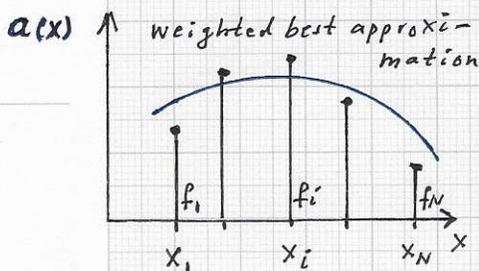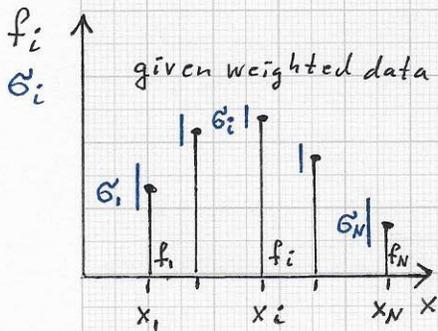• maximal difference: $\max_{\Omega} | p(x) - \hat{p}(x) |$

• discrete RMS: $\left( \frac{1}{N} \sum_{i=0}^{N} \left( p(x_i) - \hat{p}(x_i) \right)^2 \right)^{1/2}$,

$\left( \frac{1}{N} \sum_{j=1}^{\hat{N}} \left( p(\hat{x}_j) - \hat{p}(\hat{x}_j) \right)^2 \right)^{1/2}$.

# ■ GENERATION OF SYNTHETIC DATA — Cont'd.

- **Chi-square / weighted best approximation:**

The statistical properties variance ($\sigma^2$) and standard deviation ($\sqrt{\sigma^2}$) are often considered as "WEIGHTS" in the context of constructing a best approximation for a given distribution of data (given either as discrete / binned data or as a distribution function). If it is possible to determine meaningful variance / standard deviation values to discrete / binned data, or analytically defined variance / standard deviation functions to distribution functions, then one can perform WEIGHTED BEST APPROXIMATION (or CHI-SQUARE APPROXIMATION). Often the value of $\chi^2$ is minimized when performing 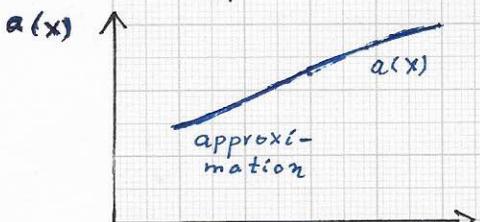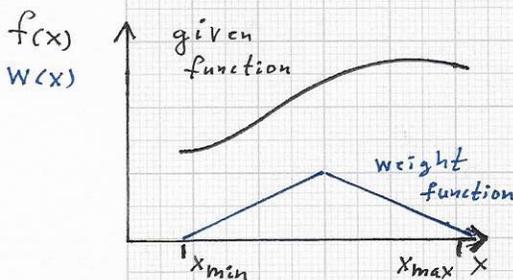best approximation of discrete, binned distribution data. Specifically, and considering a more general discrete data approximation setting, one is given values $x_i$ in the domain, dependent values $f_i$ at the locations $x_i$, and standard deviations $\sigma_i$ associated only with the $f_i$ values. Chi-square approximation minimizes the quantity

$$\chi^2 = \sum_{i=1}^{N} \left( \left( a(x_i) - f_i \right) / \sigma_i \right)^2 .$$

The function $a(x)$ is the best approximation to be computed; $1/\sigma_i^2$ serves as a "weight" for $f_i$.

given weighted data

$f_i$
$\sigma_i$

weighted best approximation

given discrete data $(x_i, f_i)$, with standard deviations $\sigma_i$ (="weights") for each $f_i$ value; weighted best approximation emphasizing low-$\sigma$ data

$f(x)$
$w(x)$

given function

weight function

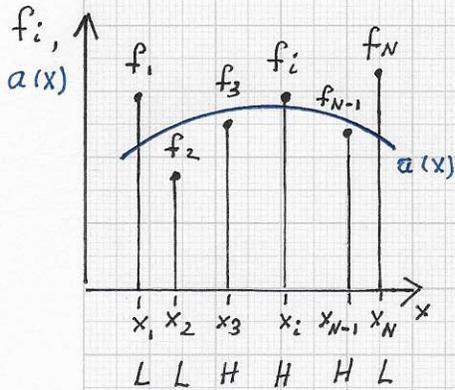$x_{min}$     $x_{max}$

$a(x)$

$a(x)$

approximation

given function and weight function; weighted best approximation emphasizing "middle part" of $f(x)$

■ <u>GENERATION OF SYNTHETIC DATA</u> — Cont'd.

• <u>Weighted best approximation</u>:

$f_i,$ $a(x)$



L L H H H L

data $(x_i, f_i)$ with
Low (L) or high (H)
weights to be appro-
ximated with $a(x)$
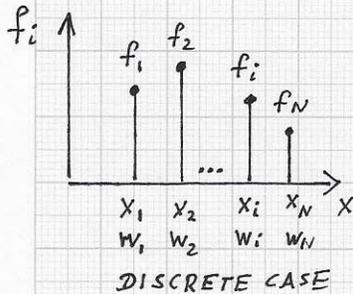— "DISCRETE CASE"

$f(x),$
$w(x),$
$a(x)$



function $f(x)$ and
weight function $a(x)$
given, weighted best
approximation $a(x)$
computed
— "CONTINUOUS CASE"

We briefly review the method of <u>weighted</u>
<u>best approximation</u>. It is important when
constructing an analytical representation
of a given distribution of feature
point/vector data via an approximation
over a region in D-dimensional feature space.
In some cases it is possible to assign <u>variance/</u>
<u>standard deviation</u> values to an <u>individual</u>
feature point $\underline{X}_i = (x_1^i, ..., x_D^i)^T$ and/or to
the associated <u>probability</u> values $p_j$ that
characterize the distribution implied by $\{\underline{X}_i\}_{i=1}^N$
via a set of $p_j$-values at certain locations $\underline{X}_j$.
The $p_j$-values must be approximated with
weights $w_j$. The weights $w_j$ are inversely
proportional to variance/standard deviation:
A datum with <u>low variance/standard deviation</u>
has <u>high weight</u> and vice versa. Thus, we
can use <u>weighted best approximation</u> to compute
the needed <u>distribution</u> function $p(\underline{X})$ from $p_j$-
and $w_j$-values given at points $\underline{X}_j$. For example,
one could use a polynomial (or piecewise polynomial)
function as approximation.
Weighted best approximation can also be used
when the given <u>data are a function $f(\underline{X})$ (to be</u>
approximated) with an associated <u>weight function</u> $W(\underline{X})$.
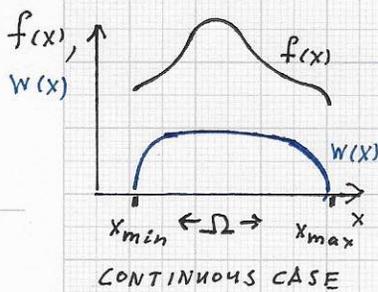The approximation $a(\underline{X})$ will "emphasize" high-weight regions.

■ GENERATION OF SYNTHETIC DATA – Cont'd.
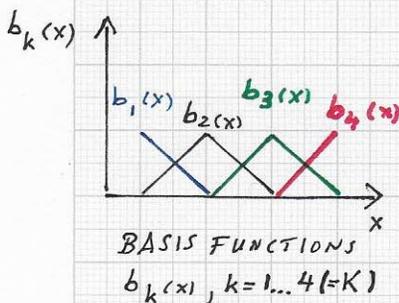
• Weighted best approximation:

One usually employs least-squares approximation to compute the unique weighted best approximation in the discrete and continuous cases. Thus, one minimizes the following quantities:
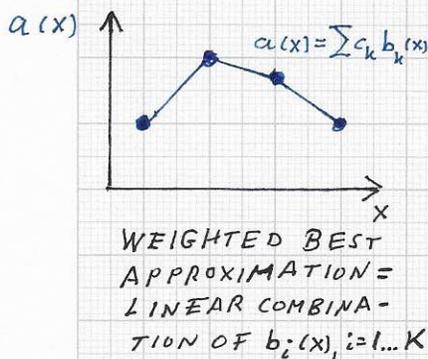
$$\left( \sum_{i=1}^{N} w_i \left( a(x_i) - f_i \right)^2 \right), \quad w_i > 0,$$

$f_i$

DISCRETE CASE

$$\left( \int_{\Omega} w(x) \left( a(x) - f(x) \right)^2 \right), \quad w(x) > 0.$$

$f(x), w(x)$

CONTINUOUS CASE

The weights $w_i > 0$ (or the weight function $w(x)$) should be "normalized" in such a way that a w-value represents the probability of a function to have the value $f_i$ ($f(x)$) at location $x_i$ ($x$).

$b_k(x)$

$b_1(x) \quad b_2(x) \quad b_3(x) \quad b_4(x)$

BASIS FUNCTIONS $b_k(x), k=1...4 (=K)$

The weighted best approximation $a(x)$ itself is a linear combination of appropriate basis functions $b_k(x)$ (e.g., polynomials, spline basis functions or radial basis functions):

$a(x)$

$a(x) = \sum c_k b_k(x)$

$$a(x) = \sum_{k=1}^{K} c_k b_k(x) .$$

WEIGHTED BEST APPROXIMATION = LINEAR COMBINATION OF $b_i(x), i=1...K$

The unknown coefficients $c_k$ result from solving the normal equations defined by inserting $\sum_{k=1}^{K} c_k b_k(x)$ into the expressions for the discrete or continuous case to be minimized.

~ BH