# Package 'primaryTranscriptAnnotation'

June 3, 2019

**Type** Package

**Title** Data-driven gene coordinate inference and assignment of
annotations to transcriptional units identified de novo

**Version** 0.1.0

**Author** Warren Anderson, Mete Civelek, Michael Guertin

**Maintainer** Warren Anderson <warrena@virginia.edu>

**Description** This package was developed for two purposes: (1) to generate data-driven transcript annotations based on nascent transcript sequencing data, and (2) to annotate de novo identified transcriptional units (e.g., genes and enhancers) based on existing annotations such as those from (1). The code for this package was employed in analyses underlying our manuscript in preparation: Transcriptional mechanisms and gene regulatory networks underlying early adipogenesis (this note will be updated upon publication).

**License** NA

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 6.1.1

**Imports** bigWig, pracma, dplyr

## R topics documented:

---

adjacent.gene.tss           *Function to identify transcription start sites (TSSs) for adjacent gene pairs*

---

### Description

We have identified adjacent gene pairs with overlaps based on manual analyses. We empirically define TSSs for these genes by binning the region spanned by both genes, fitting smooth spline curves to the binned read counts, and identifying the two largest peaks separated by a specified distance. We set a bin size and apply the constraint that the identified 'pause' peaks must be a given distance apart. The search region includes a specified distance upstream of the upstream-most gene. We shift the identified peaks upstream. For the spline fits, we set the number of knots to the number of bins divided by specified setting. See also get.TSS().

### Usage

```
adjacent.gene.tss(fix.genes = NULL, bed.long = NULL, bw.plus = NULL,
  bw.minus = NULL, bp.bin = NULL, shift.up = NULL,
  delta.tss = NULL, knot.div = 4, knot.thresh = 5, diff.tss = 1000,
  fname = "adjacentTSS.pdf")
```

### Arguments

| | |
|---|---|
| fix.genes | frame with upstream and downstream genes |
| bed.long | long gene annotations, see get.largest.interval() |
| bw.plus | plus strand bigWig data |
| bw.minus | minus strand bigWig data |
| bp.bin | the interval will be separated into adjacent bins of this size |
| shift.up | look upstream of the long gene start by this amount to search for a viable TSS |
| delta.tss | amount by which to shift the identied TSS upstream of the identified interval |
| knot.div | the number of knots for the spline fit is defined as the number of bins covering the gene end divided by this number |
| knot.thresh | minimum number of knots |
| diff.tss | minimum distance between TSSs |
| fname | file name (.pdf) for output plots |

## Value

A vector of TSSs, along with a plot. The titles for the plots indicate the upstream gene, the downstream gene, and the strand. For genes on the plus strand, the upstream gene is on the left. For the minus strand, the upstream gene is on the right. The red line denotes the spline fit and the vertical lines indicate the TSSs. The plot is intended for diagnostic purposes.

## Examples

```
# get the start sites for adjacent gene pairs
bp.bin = 10
knot.div = 40
shift.up = 100
delta.tss = 50
diff.tss = 2000
TSS.adjacent = adjacent.gene.tss(fix.genes=fix.genes, bed.long=bed.long,
                          bw.plus=bw.plus, bw.minus=bw.minus,
                          knot.div=knot.div, bp.bin=bp.bin,
                          shift.up=shift.up, delta.tss=delta.tss,
                          diff.tss=diff.tss, fname="adjacentTSS.pdf")
```

---

adjacent.gene.tts    *Function to identify gene ends (TTSs) for adjacent gene pairs*

---

## Description

We set the TTSs for the upstream genes to a given distance from the starts of the downstream genes for the adjacent gene pairs. See also adjacent.gene.tss().

## Usage

```
adjacent.gene.tts(fix.genes = NULL, bed.long = NULL,
  TSS.adjacent = NULL, dist.from.start = NULL)
```

## Arguments

| | |
|---|---|
| fix.genes | frame with upstream and downstream genes for adjacent gene pairs |
| bed.long | comprehensive gene annotations with large intervals defined by get.largest.interval() |
| TSS.adjacent | TSSs for adjacent genes, identified from adjacent.gene.tss() |
| dist.from.start | |
| | distance between the TTS of the upstream gene and the TSS of the downstream gene |

## Value

A bed6 frame with TTSs for upstream genes in adjacent gene pairs

## Examples

```
# get the end sites for adjacent upstream genes, integrate with TSSs
dist.from.start = 50
adjacent.tts.up = adjacent.gene.tts(fix.genes=fix.genes, bed.long=bed.long,
                            TSS.adjacent=TSS.adjacent, dist.from.start=dist.from.start)
tss = TSS.adjacent$TSS
names(tss) = TSS.adjacent$gene
adjacent.tts.up = apply.TSS.coords(bed=adjacent.tts.up, tss=tss)
```

---

apply.TSS.coords                    *Apply previously identified TSSs to an existing gene annotation*

---

## Description

This function will apply transcription start site (TSS) coordinates to a bed6 frame to incorporate the output from get.TSS() into an existing bed frame.

## Usage

```
apply.TSS.coords(bed = NULL, tss = NULL)
```

## Arguments

| | |
|---|---|
| bed | bed6 file with comprehensive gene annotations |
| tss | TSS estimates, output from get.TSS(). Note that the TSS genes must be a subset of the genes in the bed6 data. |

## Value

A bed6 file with estimated TSSs incorporated

## Examples

```
# apply TSS coordinates to the expression-filtered long gene annotations
bed.long.filtered.tss = apply.TSS.coords(bed=bed.long.filtered, tss=TSS.gene)
```

---

gene.end.plot *Plot the results of transcription terminination site identification*

---

## Description

This function can be used to evaluate the performance of get.gene.end() and get.TTS(). We also recommend visualizing the results of data-driven gene annotations using a genome browser.

## Usage

```
gene.end.plot(bed = NULL, gene = NULL, bw.plus = NULL,
  bw.minus = NULL, bp.bin = NULL, pk.thresh = 0.1, knot.div = 40,
  knot.thresh = 5, cnt.thresh = 5, add.to.end = 50000,
  tau.dist = 10000, frac.max = 1, frac.min = 0.2)
```

## Arguments

| | |
|---|---|
| bed | bed6 gene coordinates with gene end intervals for estimation of the TTS |
| gene | a specific gene for plotting |
| bw.plus | plus strand bigWig data |
| bw.minus | minus strand bigWig data |
| bp.bin | the interval at the gene end will be separated into adjacent bins of this size |
| pk.thresh | the TTS is defined as pk.thresh percent of max peak of the spline fit |
| knot.div | the number of knots for the spline fit is defined as the number of bins covering the gene end divided by this number; increasing this parameter results and a smoother curve |
| knot.thresh | minimum number of knots, if knot.div gives a lower number, use this |
| cnt.thresh | read count number below which the TTS is not evaluated |
| add.to.end | the maximal length of the search region (bp) |
| tau.dist | distance constant for the exponential defining the region for peak detection |
| frac.max | maximal fraction of gene end region for peal detection |
| frac.min | minimal fraction of gene end region for peal detection |

## Value

A single plot is generated

## Examples

```
gene.end.plot(bed=bed.for.tts.eval, gene="Aamp",
              bw.plus=bw.plus, bw.minus=bw.minus,
              bp.bin=bp.bin, add.to.end=add.to.end, knot.div=knot.div,
              pk.thresh=pk.thresh, knot.thresh=knot.thresh,
              cnt.thresh=cnt.thresh, tau.dist=tau.dist,
              frac.max=frac.max, frac.min=frac.min)
```

---

gene.overlaps                  *Documentation of gene overlaps*

---

### Description

This function identifies gene overlaps that do not involve identical coordinates.

### Usage

```
gene.overlaps(bed = NULL)
```

### Arguments

bed                  bed6 file with processed gene annotations

### Value

A list with has.start.inside and is.a.start.inside. has.start.inside = bed6 file that documents genes in which there are starts of other genes. is.a.start.inside = bed6 file that documents genes that start inside other genes.

### Examples

```
# run overlap analysis
overlap.data = gene.overlaps( bed = bed.long.filtered2.tss )
has.start.inside = overlap.data$has.start.inside
is.a.start.inside = overlap.data$is.a.start.inside
dim(has.start.inside)
dim(is.a.start.inside)
```

---

get.end.intervals              *Get intervals at the gene ends for transcription termination site (TTS) estimation*

---

### Description

We define regions at the gene ends to examine read counts for TTS identification. Note that transcription frequently extends beyond the poly-A site of a gene. To capture the end of transcription, it is critical to examine regions beyond annotated gene boundries. We evaluate evidence of transcriptional termination in regions extending from a 3' subset of the gene to a selected number of base pairs downstream of the most distal annotated gene end. We initiate the search region with the end of the largest gene annotation (see get.largest.interval()). We extend the search region up to a given distance past the annotated gene end. We also apply the constraint that a TTS cannot be identified closer than a specified distance to the previously identified TSS of a downstream gene. This analysis also incorporates the constraint that a gene end region identified cannot cross the TSS of a downstream gene, thereby preventing gene overlaps on a given strand. Thus, we clip the amount of bases added on to the gene end as necessary to avoid overlaps.

## Usage

```
get.end.intervals(bed = NULL, add.to.end = NULL, fraction.end = NULL,
  dist.from.start = NULL)
```

## Arguments

| | |
|---|---|
| bed | processed bed6 file with one entry per gene and identified TSSs |
| add.to.end | number of bases to add to the end of each gene to define the search region |
| fraction.end | fraction of the gene annotation to consider at the end of the gene |
| dist.from.start | |
| | the maximal allowable distance between the end of one gene and the start of the next |

## Value

A bed6 for gene end evaluation

## Examples

```
# get intervals for TTS evaluation
add.to.end = 100000
fraction.end=0.1
dist.from.start=50
bed.for.tts.eval = get.end.intervals(bed=bed.long.filtered4.tss,
                            add.to.end=add.to.end,
                            fraction.end=fraction.end,
                            dist.from.start=dist.from.start)
```

---

get.gene.end                 *Defining gene ends*

---

## Description

Given regions within which to search for TTSs, we opperationally define the TTSs by binning the gene end regions, counting reads within the bins, fitting smooth spline curves to the bin counts, and detecting points at which the curves decay towards zero. We applied the constraint that there must be a specified number of bases in the gene end interval, otherwise the TTS analysis is not applied. Similarly, if the number of knots identified is too low, then we set the number of knots to a specified threshold. For this analysis, we set a sub-region at the beginning of the gene end region and identify the maximal peak from the spline fit. Then we identify the point at which the spline fit decays to a threshold level of of the peak level. We reasoned that the sub-region should be largest for genes with the greatest numbers of clipped bases, because such cases occur when the conventional gene ends are proximal to identified TSSs, and we should include these entire regions for analysis of the TTS. Similarly, we reasoned that for genes with substantially less clipped bases, and correspondingly larger gene end regions with greater potential for observing enhancers or divergent transcripts, the sub-regions should be smaller sections of the upstream-most gene end region. We use an exponential model to define the sub-regions. The user should use the output of

get.end.intervals() as an input to this function. Note that gene ends will be set to zeros if there are
no counts or if the interval is too small. This function calls get.TTS().

## Usage

```
get.gene.end(bed = NULL, bw.plus = NULL, bw.minus = NULL,
  bp.bin = NULL, pk.thresh = 0.1, knot.div = 40, knot.thresh = 5,
  cnt.thresh = 5, add.to.end = 50000, tau.dist = 10000,
  frac.max = 1, frac.min = 0.2)
```

## Arguments

| | |
|---|---|
| bed | bed6 gene coordinates with gene end intervals for estimation of the TTS |
| bw.plus | plus strand bigWig data |
| bw.minus | minus strand bigWig data |
| bp.bin | the interval at the gene end will be separated into adjacent bins of this size |
| pk.thresh | the TTS is defined as pk.thresh percent of max peak of the spline fit |
| knot.div | the number of knots for the spline fit is defined as the number of bins covering the gene end divided by this number; increasing this parameter results and a smoother curve |
| knot.thresh | minimum number of knots, if knot.div gives a lower number, use this |
| cnt.thresh | read count number below which the TTS is not evaluated |
| add.to.end | the maximal length of the search region (bp) |
| tau.dist | distance constant for the exponential defining the region for peak detection |
| frac.max | maximal fraction of gene end region for peal detection |
| frac.min | minimal fraction of gene end region for peal detection |

## Value

A bed6 file with TTS estimates incorporated, zeros are present if the TTS could not be estimated.
The output is a list including the bed frame along with several metrics (see example below and
vignette).

## Examples

```
# identify gene ends
add.to.end = max(bed.for.tts.eval$xy)
knot.div = 40
pk.thresh = 0.02
bp.bin = 50
knot.thresh = 5
cnt.thresh = 5
tau.dist = 10000
frac.max = 1
frac.min = 0.2
gene.ends = get.gene.end(bed=bed.for.tts.eval, bw.plus=bw.plus, bw.minus=bw.minus,
                    bp.bin=bp.bin, add.to.end=add.to.end, knot.div=knot.div,
```

```
                          pk.thresh=pk.thresh, knot.thresh=knot.thresh,
                          cnt.thresh=cnt.thresh, tau.dist=tau.dist,

# get metrics
minus.lowcount = length(gene.ends$minus.lowcount)
plus.lowcount = length(gene.ends$plus.lowcount)
minus.knotmod = length(gene.ends$minus.knotmod)
plus.knotmod = length(gene.ends$plus.knotmod)
ends = gene.ends$bed
```

---

| get.largest.interval | *Get the largest interval for each gene, given multiple TSS and TTS annotations* |
| --- | --- |

---

### Description

The input bed6 file can be derived from a gencode annotation file, as described in the vignette

### Usage

```
get.largest.interval(bed = NULL)
```

### Arguments

bed                bed6 frame with comprehensive gene annotations, defaults to NULL

### Value

a bed6 frame with the largest annotation for each gene

### Examples

```
# get intervals for furthest TSS and TTS +/- interval
bed.long = get.largest.interval(bed=dat0)
```

---

| get.TSS | *Select an optimal transcription start site (TSS) for each gene* |
| --- | --- |

---

### Description

We set a range around each annotated gene start within which to search for a region of peak read density. Such regions of peak read density occur at 'pause sites' that are typically 20-80 bp from the TSS. Within each region around an annotated gene start, we translate a sliding window throughout to find the sub-region with maximal density. We determine the window with the maximal density out of all regions considered and define the TSS as the upstream boundry of this window, minus a shift of a specified amount. The output of this function should be processed using the function apply.TSS.coords().

**Usage**

```
get.TSS(bed = NULL, bw.plus = NULL, bw.minus = NULL,
  bp.range = NULL, bp.delta = NULL, bp.bin = NULL,
  delta.tss = NULL, cnt.thresh = NULL)
```

**Arguments**

| | |
|---|---|
| bed | bed6 file with comprehensive gene annotations |
| bw.plus | plus strand bigWig data |
| bw.minus | minus strand bigWig data |
| bp.range | range to be centered on each gene start for evaluating read counts. Note that the value should be divisible by 2 |
| bp.delta | number of reads to move incrementally - sliding window interval |
| bp.bin | bin size for evaluating read counts within a given region |
| delta.tss | amount by which to shift the identied TSS upstream of the identified interval |
| cnt.thresh | read count number below which the TTS is not evaluated |

**Value**

A vector of indentified transcription start site coordinates. Note that strand and chromosome information are not provided here.

**Examples**

```
# select the TSS for each gene
bp.range = 1200
bp.delta = 10
bp.bin = 50
delta.tss = 50
cnt.thresh = 5
TSS.gene = get.TSS(bed=dat0, bw.plus=bw.plus, bw.minus=bw.minus,
                bp.range=bp.range, bp.delta=bp.delta, bp.bin=bp.bin,
                delta.tss=delta.tss, cnt.thresh=cnt.thresh)
```

---

get.TTS *Estimate the transcription termination sites (TTSs)*

---

**Description**

This function estimates the TTSs, this function is called by get.gene.end() and is not designed to be run on its own.

## Usage

```
get.TTS(bed = NULL, bw = NULL, bp.bin = NULL, pk.thresh = NULL,
    knot.div = NULL, cnt.thresh = NULL, knot.thresh = NULL,
    add.to.end = NULL, tau.dist = NULL, frac.max = NULL,
    frac.min = NULL)
```

## Arguments

| | |
|---|---|
| bed | bed6 gene coordinates with gene end intervals for estimation of the TTS |
| bw | plus or minus strand bigWig data |
| bp.bin | the interval at the gene end will be separated into adjacent bins of this size |
| pk.thresh | the TTS is defined as pk.thresh percent of max peak of the spline fit |
| knot.div | the number of knots for the spline fit is defined as the number of bins covering the gene end divided by this number |
| cnt.thresh | read count number below which the TTS is not evaluated |
| knot.thresh | minimum number of knots, if knot.div gives a lower number, use this |
| add.to.end | the maximal length of the search region (bp) |
| tau.dist | distance constant for the exponential defining the region for peak detection |
| frac.max | maximal fraction of gene end region for peal detection |
| frac.min | minimal fraction of gene end region for peal detection |

## Value

A vector of TTS coordinates

## Examples

```
See get.gene.end()
```

---

| multi.overlap.assign | *Function to assign identifiers for class5 overlap TUs* |
|---|---|

---

## Description

This function is called inside of multiple.overlaps() and is not recommended to be used on its own.

## Usage

```
multi.overlap.assign(fr = NULL, tss.thresh = NULL, delta.tss = NULL,
    delta.tts = NULL, cl = NULL)
```

**Arguments**

| | |
|---|---|
| `fr` | = frame with full bedtools overlap data for a class5 TU with multiple internal annotations |
| `tss.thresh` | = number of bp a TU beginning can be off from an annotation in order to be assigned that annotation |
| `delta.tss` | = max distance between an upstream gene end and downstream gene start cl = multi-overlap class |
| `delta.tts` | max difference distance between and annotated gene end and the start of a down-stream gene before an intermediate TU id is assigned |
| `cl` | multi-overlap class |

**Value**

A bed data with chr, start, end, id, class, and strand. Note that id = 0 for regions of large TUs that do no match input annotations.

**Examples**

```
See multiple.overlaps()
```

---

   multi.overlaps                *Assign identifiers to TUs with multiple gene overlaps*

---

**Description**

This function will assign identifiers to TUs that overlapped with multiple genes. Trusted gene annotations are considered in terms of their degree of overlap with TUs identified in an unbiased manner. Marginal regions of TUs outside of gene overlaps are given generic identifiers.

**Usage**

```
multi.overlaps(overlaps = NULL, tss.thresh = NULL, delta.tss = NULL,
  delta.tts = NULL)
```

**Arguments**

| | |
|---|---|
| `overlaps` | frame with bed intersect of inferred TUs (col1-6) with annotated TUs (col7-12) and overlaping bps (col14) |
| `tss.thresh` | number of bp a TU beginning can be off from an annotation in order to be assigned that annotation |
| `delta.tss` | max distance between an upstream gene end and downstream gene start |
| `delta.tts` | max difference distance between and annotated gene end and the start of a down-stream gene before an intermediate TU id is assigned |

## Value

A list with bed data (bed) and counts for each class (cnt5-8). bed = data with chr, start, end, id, class, and strand. cnt5 = count of TU from class 5 overlaps

## Examples

```
dup.hmm.rows = overlap.tu[duplicated(overlap.tu[,c(1:3,6)]) |
    duplicated(overlap.tu[,c(1:3,6)], fromLast=TRUE),]
dup.full = dup.hmm.rows
names(dup.full)[1:6] = c("infr.chr","infr.start","infr.end",
                         "infr.gene","infr.xy","infr.strand")
nrow(dup.full[,1:3] %>% unique)
nrow(dup.full %>% select(ann.gene) %>% unique)
tss.thresh = 200
delta.tss = 50
delta.tts = 1000
class5678 = multi.overlaps(overlaps=dup.full, tss.thresh=tss.thresh,
                           delta.tss=delta.tss, delta.tts=delta.tts)
new.ann.mult = class5678$bed
```

---

read.count.end          *Select a regions containing the gene ends*

---

## Description

Select a fraction of the annotated gene end and consider an additional number of base pairs beyond the gene end within which to count reads

## Usage

```
read.count.end(bed = NULL, bw.plus = NULL, bw.minus = NULL,
  fraction.end = NULL, add.to.end = NULL)
```

## Arguments

| | |
|---|---|
| bed | a bed6 frame |
| bw.plus | plus strand bigWig data |
| bw.minus | minus strand bigWig data |
| fraction.end | fraction of the annotated gene end (0,1) |
| add.to.end | number of base pairs beyond the gene end within which to count reads |

## Value

a list with counts and desity. counts = vector with end of gene counts for each gene. density = vector with end of gene densities for each gene.

## Examples

```
# get read counts and densities at the end of each annotated gene
fraction.end = 0.1
add.to.end = 0
end.reads = read.count.end(bed=bed.long, bw.plus=bw.plus, bw.minus=bw.minus,
                           fraction.end=fraction.end, add.to.end=add.to.end)
hist(log(end.reads$density))
hist(log(end.reads$counts))
```

---

single.overlap.assign   *Function to assign identifiers for class1 overlap TUs*

---

## Description

This function is called inside of single.overlaps() and is not recommended to be used on its own.

## Usage

```
single.overlap.assign(fr = NULL, tss.thresh = NULL, delta.tss = NULL,
  delta.tts = NULL, cl = NULL)
```

## Arguments

| | |
|---|---|
| fr | frame with full bedtools overlap data for a class1 TU with multiple internal annotations |
| tss.thresh | number of bp a TU beginning can be off from an annotation in order to be assigned that annotation |
| delta.tss | max distance between an upstream gene end and downstream gene start |
| cl | multi-overlap class |

## Value

A bed data frame with chr, start, end, id, class, and strand. Note that id = 0 for regions of large TUs that do no match input annotations.

## Examples

```
See single.overlaps()
```

---

single.overlaps                *Assign identifiers to TUs with single gene overlaps*

---

### Description

This function will assign identifiers to TUs that overlapped with single genes. Trusted gene annotations are considered in terms of their degree of overlap with TUs identified in an unbiased manner. Marginal regions of TUs outside of gene overlaps are given generic identifiers.

### Usage

```
single.overlaps(overlaps = NULL, tss.thresh = NULL, delta.tss = NULL,
  delta.tts = NULL)
```

### Arguments

| | |
|---|---|
| overlaps | frame with bed intersect of inferred TUs (col1-6) with annotated TUs (col7-12) and overlaping bps (col14) |
| tss.thresh | = number of bp a TU beginning can be off from an annotation in order to be assigned that annotation |
| delta.tss | max distance between an upstream gene end and downstream gene start |
| delta.tts | max difference distance between and annotated gene end and the start of a downstream gene before an intermediate TU id is assigned |

### Value

A list with bed data (bed) and counts for each class (cnt1-4). bed = data with chr, start, end, id, class, and strand. cnt1 = count of TU from class 1 overlaps

### Examples

```
sing.hmm.rows = overlap.tu[!duplicated(overlap.tu[,c(1:3,6)]) &
     !duplicated(overlap.tu[,c(1:3,6)], fromLast=TRUE),]
overlaps = sing.hmm.rows
names(overlaps)[1:6] = c("infr.chr","infr.start","infr.end",
                         "infr.gene","infr.xy","infr.strand")
tss.thresh = 200
delta.tss = 50
delta.tts = 1000
class1234 = single.overlaps(overlaps=overlaps, tss.thresh=tss.thresh,
                         delta.tss=delta.tss, delta.tts=delta.tts)
new.ann.sing = class1234$bed
```

---

TSS.count.dist                      *Get distances from identified transcription start sites (TSSs) to the*
                                    *nearest regions with peak reads*

---

**Description**

This function was designed to evaluate the results of out TSS identification analysis. We specify a
region centered on the TSS. We obtain read counts in bins that span this window. We sort the genes
based on the bin with the maximal reads and we scale the data to the interval (0,1) for visualization.
We compute the distances between the TSS and the bin with the max reads within the specified
window.

**Usage**

```
TSS.count.dist(bed = NULL, bw.plus = NULL, bw.minus = NULL,
  window = NULL, bp.bin = NULL)
```

**Arguments**

| | |
|---|---|
| bed | = bed file for TSS evaluation |
| bw.plus | = plus strand bigWig data |
| bw.minus | = minus strand bigWig data |
| window | = region size, centered on the TSS for analysis |
| bp.bin | = the interval will be separated into adjacent bins of this size |

**Value**

A list with three vector elements: raw, scaled, and dist. All outputs are organized based on TSS
position (left-upstream). raw = binned raw counts. scaled = binned scaled counts (0,1). dist = upper
bound distances ( min(|dists|) = bp.bin ). See examples and vignette for more details.

**Examples**

```
# look at read distribution around identified TSSs
bp.bin = 10
window = 1000
tss.dists = TSS.count.dist(bed=bed.tss.tts, bw.plus=bw.plus, bw.minus=bw.minus,
                           window=window, bp.bin=bp.bin)

# look at read distribution around 'long gene' annotation TSSs
tss.dists.lng = TSS.count.dist(bed=bed.long[bed.long$gene %in% names(tss.dists$dist),],
                               bw.plus=bw.plus, bw.minus=bw.minus,
                               window=window, bp.bin=bp.bin)
```

---

`TTS.count.dist`  *Get read counts around transcription termination sites (TTSs)*

---

### Description

We set a window around the TTS and segment the window into bins. To sort the read data, we compute cumulative counts of reads, and we sort based on the bin at which a specified percentage of the reads are found. We also take a ratio of gene counts downstream / upstream of the TTS with regions within an interval.

### Usage

```
TTS.count.dist(bed = NULL, bw.plus = NULL, bw.minus = NULL,
  window = NULL, bp.bin = NULL, frac.max = NULL,
  ratio.region = NULL)
```

### Arguments

| | |
|---|---|
| bed | bed frame for TSS evaluation |
| bw.plus | plus strand bigWig data |
| bw.minus | minus strand bigWig data |
| window | region size, centered on the TTS for analysis |
| bp.bin | the interval will be separated into adjacent bins of this size |
| frac.max | fraction of cumulative distribution for sorting entries are sorted by indices min( cumsum(x)/sum(x) ) < frac.max |
| ratio.region | number of bp on either side of the TTS to evaluate the ratio entries |

### Value

A list with three elements: raw, scaled, and ratio. All outputs are organized based on TTS position (left-upstream, see frac.max). raw = binned raw counts. scaled = binned scaled counts (0,1). ratio = downstream(ratio.region) / upstream(ratio.region).

### Examples

```
# look at read distribution around identified TTSs
window = 1000
bp.bin = 10
frac.max = 0.8
ratio.region = 300
tts.dists = TTS.count.dist(bed=bed.tss.tts, bw.plus=bw.plus, bw.minus=bw.minus,
    window=window, bp.bin=bp.bin, frac.max=frac.max,
    ratio.region=ratio.region)
```

# Index